

---

# Thermodynamics of protein folding from coarse-grained models' perspectives

Michael Bachmann and Wolfhard Janke

Institut für Theoretische Physik and Centre for Natural Sciences (NTZ),  
Universität Leipzig, Postfach 100 920, D-04009 Leipzig, Germany  
{michael.bachmann,wolfhard.janke}@itp.uni-leipzig.de

## 1 Introduction

Proteins are linear chains of amino acids connected by covalent peptide bonds. Twenty types of amino acids, mainly differing in the molecular structure of their side chains, were identified in bioproteins. Since bioproteins typically consist of hundreds to thousands of amino acids, the number of possible amino acid sequences is extremely large. Considering, for example, a chain of only 100 amino acids, the number of possible sequences (i.e., primary structures) is of order  $10^{130}$ , but this is only one side of the medal. The main importance of the proteins lies in their function in a biological organism, and the function is inseparably connected with the geometrical structure of the protein, i.e., its folded, native conformation which is usually divided into secondary, tertiary, and quarternary substructures [1,2]. The required stability of this native state against thermal and other environmental fluctuations rules out most of the possible sequences [3]. It is not yet understood, how the relatively small number of relevant proteins, e.g., about  $10^5$  in the human body, has been selected by nature in an evolutionary process [4].

The physical interactions responsible for the folding of a protein into its native structure are in principle known. The complexity of the macromolecule with up to ten thousands of atoms, however, makes precise predictions of the energetically most-favored structure based on *ab initio* quantum-mechanical calculations practically impossible. This is due to the long-range overlap of many-body electronic orbitals and the screening by the positively charged cores. The problem is indeed still more complex as the natural environment of proteins is a polar aqueous solvent. For this reason, classical models with hundreds of effective parameters ("force fields") have been developed in the past decades in order to predict native structures and to study folding dynamics in computer simulations [5]. Despite the simplifications, these models are still highly complex and hard to manage even by means of sophisticated algorithms and powerful capability computers. Furthermore it turned out that

folding and misfolding depend sensitively on the choice of the force field parameters with the consequence that predictions of different established models do frequently not coincide. Another problem is that investigations of these models require enormous computational capacities. For this reason and the fact that folding times in nature range from milliseconds to seconds, molecular dynamics simulations (MD) for studying the deterministic folding dynamics are currently widely useless as the time-scale of nano- to microseconds of reliable MD simulations is orders of magnitudes smaller.

It should be noted, however, that MD is quite successful in studies of biological short-time processes, where the biological function of proteins can be studied. Fascinating examples, where MD proved to be highly successful, are the penetration of water molecules into a cell through aquaporin being a membrane protein [6] and the ATP synthase, a process, where the catalytic subunits of F1, embedded into the membrane F0 proton channel, partially act as rotating “molecular motor” that promotes dehydration of ADP and P to ATP [7]. Such studies require that the native folds of the proteins must be known as these are used as *input*. For considering thermodynamics, Monte Carlo simulations of these all-atom models are much more promising, in particular by applying sophisticated generalized-ensemble algorithms [8]. Nonetheless, the enormous efforts required to obtain trustworthy results with these models strongly limit the systematic exposure of the general principles behind protein folding processes, which necessitates comparative studies of an appropriate set of different sequences.

In these lecture notes, we therefore follow a different approach and discuss minimalistic, coarse-grained protein models. Coarse-graining of models, i.e., increasing relevant length scales by reducing the number of microscopic degrees of freedom, has proven to be very successful in polymer science. Although specificity is much more sensitive for proteins, since details (charges, polarity, etc.) and differences of the amino acid side chains can have strong influences on the fold, mesoscopic approaches are also of essential importance for the basic understanding of conformational transitions affecting the folding process. It is also the only possible approach for systematic analyses such as the evolutionarily significant question why only a few sequences in nature are “designing”, i.e., relevant for selective functions. On the other hand, what is the reason why proteins prefer a comparative small set of target structures, i.e., what explains the preference of designing sequences to fold into the *same* three-dimensional structure? All these questions are widely still unanswered yet.

As a first step towards their solution we discuss in the first part simple hydrophobic-polar (HP) lattice models, where only the most characteristic hydrophobic or polar nature of the 20 naturally occurring amino acids is taken into account and the linear chain is modeled by a self-avoiding walk. Such models allow a comprising analysis of both, the conformation *and* sequence space, e.g., by exactly enumerating all combinatorial possibilities. Other important aspects in lattice model studies are the identification of lowest-energy

conformations of comparatively long sequences and the characterization of the folding thermodynamics.

In the second part we focus on simple AB off-lattice models, where similar to the HP model (for historical reasons)  $A$  symbolizes hydrophobic and  $B$  polar regions of the protein, whose conformations are modeled by polymer chains governed by bending energy and van der Waals interactions. These models allow for the analysis of different mutated sequences with respect to their folding characteristics. Here, the idea is that the folding transition is a kind of pseudophase transition which can in principle be described by one or a few order-like parameters. Depending on the sequence the folding process can be highly cooperative (single-exponential), less cooperative depending on the height of a free-energy barrier (two-state folding), or even frustrating due to the existence of different barriers in a metastable regime (crystal or glassy phases). These characteristics known from functional proteins can be recovered in the AB model, which is computationally much less demanding than all-atom formulations and thus enables throughout theoretical analyses.

Such coarse-grained models enable a broader view on the general problem of protein folding, but for precise, specific predictions, their applicability is limited. In analogy to magnetic systems, they are rather comparable with the Ising model for ferromagnets or the Edwards-Anderson-Ising model for spin glasses. It should also be remarked that, due to their nontrivial simplicity, coarse-grained models are also a perfect testing ground for newly developed algorithms.

## 2 Why coarse-graining?

Functional proteins in a biological organism are typically characterized by a unique three-dimensional molecular structure, which makes the protein selective for individual functions, e.g., in catalytic, enzymatic, and transport processes. In most cases, the free-energy landscape is believed to exhibit a rough shape with a large number of local minima and, for functional proteins, a deep, funnel-like global minimum. This assumed complexity is the reason, why it is difficult to understand how the random-coil conformation of covalently bonded amino acids – the sequence is generated in the ribosome according to a certain genetic sequence in the DNA – spontaneously folds into a well-defined stable “native” conformation. Furthermore, it is expected that there are only a small number of folding paths from any unfolded conformation to this final fold.

Protein folding follows a strict hierarchy at different length scales. The so-called primary structure, i.e., the sequence of amino acids in the linear chain is provided by the ribosome. Since subsequent amino acids are uniformly linked by a covalent peptide bond independent of the geometrical structure of the protein, the typical length scale of the primary structure is a single amino acid. The next level are secondary structures like  $\alpha$ -helices,  $\beta$ -sheets, and

turns. The main reason for the formation of these structures is backbone hydrogen bonding which typically involves segments of several subsequent amino acids. Therefore, the scale of secondary structures is determined by the typical segment sizes, which are of the order of ten amino acids. Consequently, secondary-structure formation is the first step in protein folding. This is followed by the formation of global, single-domain tertiary structures. In fact, this process is what renders protein folding special. The main driving force for the folding of a complex domain, i.e., of up to hundreds of amino acids, is an effective cooperative intrinsic interaction between many amino acid side-chains and which is strongly influenced by the solubility properties (in particular its polarization) of the aqueous solvent the protein resides in. Roughly, amino acid side-chains can be classified as polar, hydrophobic, and neutral. While polar residues favor contact with polar water molecules, hydrophobic acids avoid contact with water which results in an effective attraction between hydrophobic side-chains. In consequence, this attractive force leads to a formation of a highly compact hydrophobic core, which is screened from the solvent by a shell of polar amino acids. For very large proteins, the final stage in the folding process is the arrangement of several domains in a quarternary structure.

Thus, the most complex process in protein folding is the formation of tertiary hydrophobic-core structures. Although atomic details, e.g., van der Waals volume exclusion separating side-chains in linear and ring structures, polarizability, and partial charges, noticeably influence the folding process and the native fold, it should be possible to understand certain aspects of the folding characteristics, at least qualitatively, by means of coarse-grained models which are based on a few effective parameters. In the following, we investigate this question within the two minimalistic HP lattice and AB off-lattice heteropolymer models.

### 3 The hydrophobic-polar (HP) lattice protein model

The simplest model for a qualitative description of protein folding is the lattice hydrophobic-polar (HP) model [9]. In this model, the continuous conformational space is reduced to discrete regular lattices and conformations of proteins are modeled as self-avoiding walks restricted to the lattice. Assuming that the hydrophobic interaction is the most essential force towards the native fold, sequences of HP proteins consist of only two types of monomers (or classes of amino acids): Amino acids with high hydrophobicity are treated as hydrophobic monomers ( $H$ ), while the class of polar (or hydrophilic) residues is represented by polar monomers ( $P$ ). In order to achieve the formation of a hydrophobic core surrounded by a shell of polar monomers, the interaction between hydrophobic monomers is attractive and short-range. In the standard formulation of the model [9], all other interactions are neglected. Variants of

the HP model also take into account (weaker) interactions between  $H$  and  $P$  monomers as well as between polar monomers [4].

Although the HP model is extremely simple, it has been proven that identifying native conformations is an NP-complete problem in two and three dimensions [10]. Therefore, sophisticated algorithms were developed to find lowest-energy states for chains of up to 136 monomers. The methods applied are based on very different algorithms, ranging from exact enumeration in two dimensions [11,12] and three dimensions on cuboid (compact) lattices [4,13–15], and hydrophobic-core construction methods [16,17] over genetic algorithms [18–22], Monte Carlo simulations with different types of move sets [23–26], and generalized ensemble approaches [27] to Rosenbluth chain-growth methods [28] of the '*Go with the Winners*' type [29–35]. With some of these algorithms, thermodynamic quantities of lattice heteropolymers were studied as well [14,27,31,34–36].

### 3.1 The HP model

A monomer of an HP sequence  $\sigma = (\sigma_1, \sigma_2, \dots, \sigma_N)$  is characterized by its residual type ( $\sigma_i = P$  for polar and  $\sigma_i = H$  for hydrophobic residues), the position  $1 \leq i \leq N$  within the chain of length  $N$ , and the spatial position  $\mathbf{x}$  to be measured in units of the lattice spacing. A conformation is then symbolized by the vector of the coordinates of successive monomers,  $\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)$ . We denote by  $x_{ij} = |\mathbf{x}_i - \mathbf{x}_j|$  the distance between the  $i$ th and the  $j$ th monomer. The bond length between adjacent monomers in the chain is identical with the spacing of the used regular lattice with coordination number  $q$ . These covalent bonds are thus not stretchable. A monomer and its nonbonded nearest neighbors may form so-called contacts. Therefore, the maximum number of contacts of a monomer within the chain is  $(q - 2)$  and  $(q - 1)$  for the monomers at the ends of the chain. To account for the excluded volume, lattice proteins are self-avoiding, i.e., two monomers cannot occupy the same lattice site. The total energy for an HP protein reads in energy units  $\varepsilon_0$  (we set  $\varepsilon_0 = 1$  in the following)

$$E_{\text{HP}} = \varepsilon_0 \sum_{\langle i, j \rangle > i+1} C_{ij} U_{\sigma_i \sigma_j}, \quad (1)$$

where  $C_{ij} = (1 - \delta_{i+1, j})\Delta(x_{ij} - 1)$  with

$$\Delta(z) = \begin{cases} 1, & z = 0, \\ 0, & z \neq 0 \end{cases} \quad (2)$$

is a symmetric  $N \times N$  matrix called *contact map* and

$$U_{\sigma_i \sigma_j} = \begin{pmatrix} u_{HH} & u_{HP} \\ u_{HP} & u_{PP} \end{pmatrix} \quad (3)$$

is the  $2 \times 2$  interaction matrix. Its elements  $u_{\sigma_i \sigma_j}$  correspond to the energy of  $HH$ ,  $HP$ , and  $PP$  contacts. For labeling purposes we shall adopt the convention that  $\sigma_i = 0 \hat{=} P$  and  $\sigma_i = 1 \hat{=} H$ .

In the simplest formulation [9], only the attractive hydrophobic interaction is nonzero,  $u_{HH}^{\text{HP}} = -1$ , while  $u_{HP}^{\text{HP}} = u_{PP}^{\text{HP}} = 0$ . Therefore,  $U_{\sigma_i \sigma_j}^{\text{HP}} = -\delta_{\sigma_i H} \delta_{\sigma_j H}$ . This parameterization, which we will traditionally call the *HP model* in the following, has been extensively used to identify ground states of HP sequences, some of which are believed to show up qualitative properties comparable with realistic proteins whose 20-letter sequence was transcribed into the 2-letter code of the HP model [16,18,37–39].

This simple form of the standard HP model suffers, however, from the fact that the lowest-energy states are usually highly degenerate and therefore the number of designing sequences (i.e., sequences with unique ground state – up to the usual translational, rotational, and reflection symmetries) is very small, at least on the three-dimensional simple cubic (sc) lattice. Incorporating additional inter-residue interactions, symmetries are broken, degeneracies are smaller, and the number of designing sequences increases [14,15]. Based on the Miyazawa-Jernigan matrix [40] of inter-residue contact energies between real amino acids, an additional attractive nonzero energy contribution for contacts between  $H$  and  $P$  monomers is more realistic [4]. In the following, we set the elements of the interaction matrix (3) to  $u_{HH}^{\text{MHP}} = -1$ ,  $u_{HP}^{\text{MHP}} = -1/2.3 \approx -0.435$ , and  $u_{PP}^{\text{MHP}} = 0$ , corresponding to Ref. [4]. The factor 2.3 is a result of an analysis for the inter-residue energies of contacts between hydrophobic amino acids and contacts between hydrophobic and polar residues [40] which motivated the relation  $2u_{HP} > u_{PP} + u_{HH}$  [4]. In the following we call this variant the *MHP model* (mixed HP model).

### 3.2 Exact enumerations for short HP sequences

The most important advantage of lattice HP-type models compared with other, more complex protein models is that it allows for comprising analyses of conformation and sequence space. This is essential for systematic studies following two main strategies in understanding protein structure formation: *direct* and *inverse* folding. Direct folding is sequence-based, i.e., the amino acid sequence is given and the global free-energy minimum conformation(s) are sought. In the inverse folding problem, a target structure is given and the question is for how many sequences this structure is the global free-energy minimum conformation.

Since it is widely believed that for bioproteins the unique global free-energy minimum conformation under physiological conditions (i.e., the native fold) is identical with the conformation of lowest total (free) energy, it is assumed that qualitative folding-related properties of HP lattice protein sequences are comparable with realistic proteins if their groundstate is nondegenerate. An HP sequence with a unique native fold is called *designing*. On the other hand,

**Table 1.** Number of designing sequences  $S_N$  (only relevant sequences [41]) in the HP and MHP model on the simple cubic lattice.

$N$	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19
$S_N^{\text{HP}}$	3	0	0	0	2	0	0	0	2	0	1	1	1	8	29	47
$S_N^{\text{MHP}}$	7	0	0	6	13	0	11	8	124	14	66	97	486	2196	9491	4885

**Table 2.** Number of designable conformations  $D_N$  (without conformations trivially symmetric by translations, rotations, and reflections) in the HP and MHP model on the simple cubic lattice.

$N$	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19
$D_N^{\text{HP}}$	1	0	0	0	2	0	0	0	2	0	1	1	1	8	28	42
$D_N^{\text{MHP}}$	1	0	0	2	2	0	5	6	30	8	31	58	258	708	1447	1623

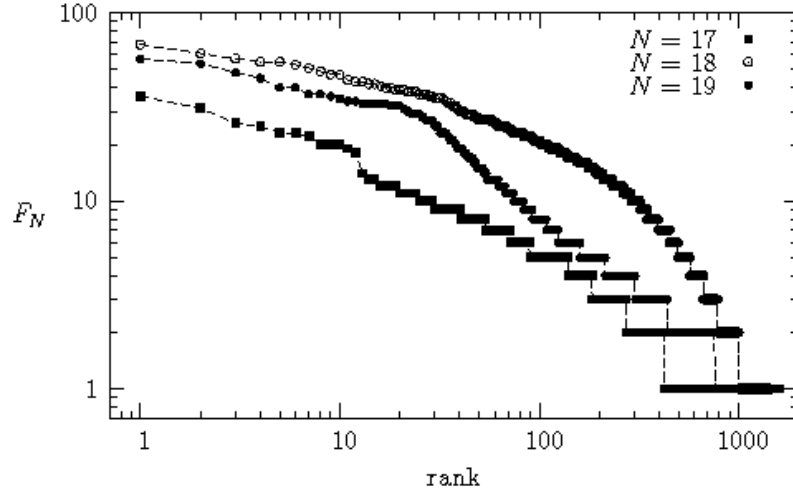
a target structure which is the native fold of one or more designing sequences, is called a *designable conformation*.

In Table 1 we list for all chain lengths  $N = 4, \dots, 19$  the total numbers  $S_N$  of relevant designing sequences [41] in the HP and the MHP model. These results were obtained by exhaustive exact enumerations of the complete conformation and sequence spaces of chains with up to 19 monomers on the sc lattice [14]. Note that there are for a 19-mer more than  $5 \times 10^5$  HP sequences and about  $2 \times 10^{12}$  self-avoiding conformations on the sc lattice, which in total allows naively more than  $10^{18}$  possible combinations. In order to achieve this, an efficient parallel implementation based on contact sets [12,42] together with symmetry considerations had to be used [15]. As already mentioned, the number of designing sequences is rather small in the standard HP model, whereas the additional *HP* attraction in the MHP model dissolves degeneracies which consequently entails a noticeably larger number of sequences with a unique ground-state conformation.

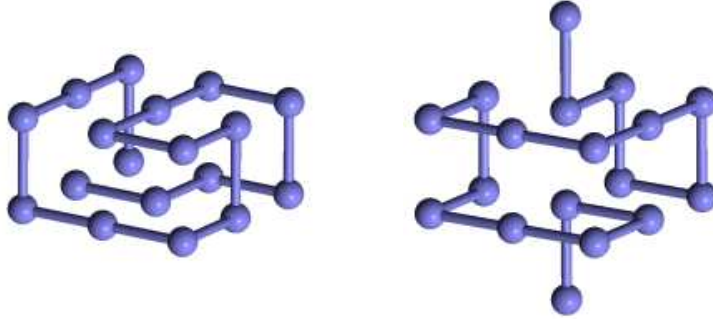
In Table 2 we list for both models the number of *different* native conformations  $D_N$  on the sc lattice. Interestingly, this number is usually much smaller than the number of designing sequences in Table 1, i.e., several designing sequences share the same ground-state conformation. The number of designing sequences that fold into a certain given target conformation  $\mathbf{X}^{(0)}$  (or conformations being trivially symmetric to this by translations, rotations, and reflections) is called *designability* [43]:

$$F_N(\mathbf{X}^{(0)}) = \sum_{\sigma \in \mathbf{S}_N} \Delta(\mathbf{X}_{\text{gs}}(\sigma) - \mathbf{X}^{(0)}), \quad (4)$$

where  $\mathbf{X}_{\text{gs}}(\sigma)$  is the native (ground-state) conformation of a designing sequence  $\sigma$  in the set of all designing sequences  $\mathbf{S}_N$  of length  $N$ . The function  $\Delta(\mathbf{Z})$  is the generalization of Eq. (2) to  $3N$ -dimensional vectors. It is unity for  $\mathbf{Z} = \mathbf{0}$  and zero otherwise.



**Fig. 1.** Designability  $F_N$  of native conformations in the MHP model for  $N = 17, 18$ , and  $19$ . The abscissa is the rank obtained by ordering all designable conformations according to their designability.



**Fig. 2.** Structure ( $N = 18$ ) with the highest designability of all native conformations (left) and most compact structure with minimal radius of gyration (right).

The designability is plotted in Fig. 1 for all native conformations that HP proteins with  $N = 17, 18$ , and  $19$  monomers can form in the MHP model. In this figure, the abscissa is the rank of the conformations, ordered according to their designability. The conformation with the lowest rank is therefore the most designable structure and we see that most of the designing sequences fold into a few number of highly designable conformations, while only a small number of designing sequences possesses a native conformation with low designability (note that the plot is logarithmic). Similar results were found, for example in Ref. [44], where the designability of compact conformations on cuboid lattices was investigated in detail. The left picture in Fig. 2 shows the



conformation with the lowest rank (or highest designability) with  $N = 18$  monomers. Note that this is not the most compact structure, i.e., the conformation with minimal gyration radius, which is shown for  $N = 18$  in Fig. 2 (right).

### 3.3 Chain-growth methods for long HP sequences

Combined exact enumeration studies of conformation *and* sequence space for lattice peptides noticeably longer than 19 monomers are currently computationally out of reach which is due to the exponential growth of the state space. Therefore, for longer sequences, primarily the direct folding problem is studied using computer simulation methods: Low-lying energetic conformations and thermodynamic properties governing the folding kinetics are identified and analyzed for a given HP sequence.

Computer simulations of lattice peptides, which are modeled as self-avoiding walks on the underlying lattice, are demanding. The reason is that the native fold, i.e., the ground-state or lowest-energy conformation, plays an essential role in protein science and that it is, in the discrete lattice representation, non- or low-degenerate. Monte Carlo simulations with move sets consisting of semilocal conformational updates like end flips, corner flips, and “crank shafts” [23–26], as well as nonlocal pivot updates [45], are inefficient in sampling the dominating dense conformations in the low-temperature region. It turned out that a different method, Rosenbluth chain growth [28] combined with a ‘Go with the winners’ strategy [29], is much more efficient in sampling highly dense conformations.

#### Pruned-enriched Rosenbluth chain-growth method (PERM)

In naive chain-growth methods based on simple sampling, a polymer grows by attaching the  $n$ th monomer at a randomly chosen nearest-neighbor site of the  $(n - 1)$ th monomer. The growth is stopped, if the total length  $N$  of the chain is reached or the randomly selected continuation of the chain is already occupied. In both cases, the next chain is started to grow from the first monomer. This simple chain growth is not yet very efficient, since the number of discarded chains grows exponentially with the chain length.

The performance can be improved with the Rosenbluth chain growth method [28], where first the free next neighbors of the  $(n - 1)$ th monomer are determined and then the new monomer is placed to one of the unoccupied sites. Since the probability for the next monomer to be set varies with the number of free neighbors, this implies a bias given by

$$p_n = \left( \prod_{l=2}^n m_l \right)^{-1}, \quad (5)$$



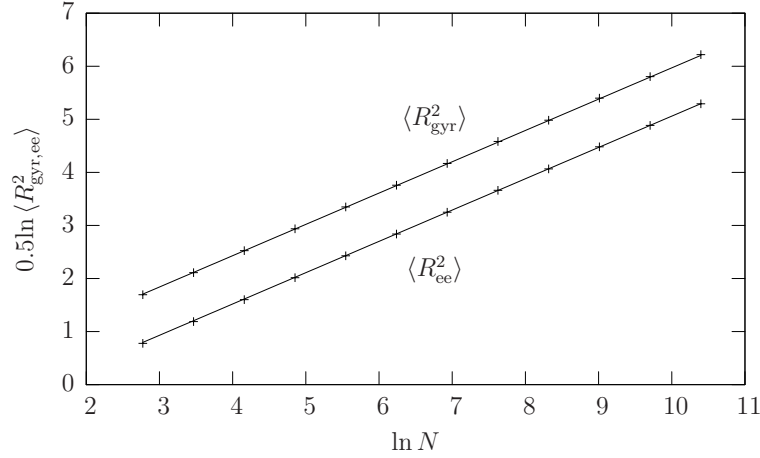
**Fig. 3.** Square-lattice example for the bias implied by Rosenbluth sampling. Both walks shown are grown from the monomer labeled “1”. Although the shapes are identical, they are created with different probabilities (left:  $p = 1/108$ , right:  $p = 1/72$ ).

where  $m_l$  is the number of free neighbors to place the  $l$ th monomer. The bias is corrected by assigning a Rosenbluth weight factor  $W_n^R = p_n^{-1}$  to each chain that has been generated by this procedure. An illustrative example for the bias in the Rosenbluth chain-growth method is shown in Fig. 3. The two depicted linear chains are grown on a square lattice from both ends (labeled by “1”). According to Rosenbluth sampling, the chain is continued if the number of free neighbor sites is  $m \geq 1$ . Since the number of free nearest-neighbor places varies, different probabilities for the continuation of the chain occur. Since both conformations are identical, the probability of creation should be the same. This requires the introduction of the correction weights. Although this biased growth is more efficient than simple sampling, this method suffers from attrition too: If all nearest neighbors are occupied, i.e., the chain was running into a “dead end” (attrition point), the complete chain has to be discarded and the growth process has to be started anew.

Combining the Rosenbluth chain growth method with population control, however, as is done in PERM (Pruned-Enriched Rosenbluth Method) [30–32], leads to a further considerable improvement of the efficiency by increasing the number of successfully generated chains. This method renders particularly useful for studying the  $\Theta$  point of polymers, since then the Rosenbluth weights of the statistically relevant chains approximately cancel against their Boltzmann probability. The (a-thermal) Rosenbluth weight factor  $W_n^R$  is therefore replaced by

$$W_n^{\text{PERM}} = \prod_{l=2}^n m_l e^{-(E_l - E_{l-1})/k_B T}, \quad 2 \leq n \leq N \quad (E_1 = 0, \quad W_1^{\text{PERM}} = 1),$$

where  $T$  is the temperature and  $E_l$  is the energy of the partial chain  $\mathbf{X}_l = (\mathbf{x}_1, \dots, \mathbf{x}_l)$  created with Rosenbluth chain growth. In PERM, population control works as follows. If a chain has reached length  $n$ , its weight  $W_n^{\text{PERM}}$  is calculated and compared with suitably chosen upper and lower threshold values,  $W_n^>$  and  $W_n^<$ , respectively. For  $W_n^{\text{PERM}} > W_n^>$ , identical copies are created which grow then independently. The weight is equally divided among them. If  $W_n^{\text{PERM}} < W_n^<$ , the chain is pruned with some proba-



**Fig. 4.** Scaling of mean square radius of gyration  $\langle R_{\text{gyr}}^2 \rangle$  and end-to-end distance  $\langle R_{\text{ee}}^2 \rangle$  for self-avoiding walks. Data points refer to results from PERM runs for  $N = 16, 32, \dots, 32768$  steps. Lines manifest the respective power-law behaviors.

bility, say  $1/2$ , and in case of survival, its weight is doubled. For a value of the weight lying between the thresholds, the chain is simply continued without enriching or pruning the sample. The upper and lower thresholds  $W_n^>$  and  $W_n^<$  are empirically parameterized. Although their values do not influence the validity of the method, a careful choice can drastically improve the efficiency of the method (the “worst” case is  $W_n^> = \infty$  and  $W_n^< = 0$ , in which case PERM is simply identical with Rosenbluth sampling). An efficient way of parameterization is dynamical adaption of the values [30–35] with respect to the actual number of generated chains  $c_n$  with length  $n$  and their estimated partition sum

$$Z_n = \frac{1}{c_1} \sum_t W_n^{\text{PERM}}(t), \quad (6)$$

where  $c_1$  is the number of growth starts (also called “tours”) and  $t$  counts the generated conformations with  $n$  monomers. Useful choices of the threshold values are

$$W_n^> = C_1 Z_n \frac{c_n^2}{c_1^3}, \quad W_n^< = C_2 W_n^>, \quad (7)$$

where  $C_1, C_2 \leq 1$  are constants. For the first tour,  $W_n^> = \infty$  and  $W_n^< = 0$ , i.e., no pruning and enriching.

In the recently developed new variants nPERMss and nPERMis [33], the number of copies is not constant and depends on the ratio of the weight  $W_n^{\text{PERM}}$  compared to the upper threshold value  $W_n^>$  and the copies are necessarily chosen to be different. The method of selecting the copies is based on simple sampling (ss) in nPERMss and a kind of importance sampling (is) in

nPERMis. This proves quite useful in producing highly compact polymers and therefore these new methods are very powerful in determining lowest-energy states of lattice proteins.

Results of a simple application of PERM to self-avoiding walks on a simple-cubic lattice are plotted in Fig. 4, where the scaling behavior  $\langle R_{\text{gyr,ee}}^2 \rangle \sim N^{2\nu}$  of the mean square radius of gyration  $\langle R_{\text{gyr}}^2 \rangle$  and end-to-end distance  $\langle R_{\text{ee}}^2 \rangle$  with the number of steps  $N$  is shown. Data were obtained for chains of  $N = 16, 32, \dots, 32768$  steps. For both quantities, the slope of the lines in the logarithmic plot is  $\nu = 0.59$ , which is close to the precisely known critical exponent  $\nu = 0.588\dots$  [47].

### Multicanonical chain-growth algorithm

The efficiency of PERM depends on the simulation temperature. Therefore, a precise estimation of the density of states requires separate simulations at different temperatures. Then, the density of states can be constructed by means of the multiple-histogram reweighting method [46]. Although being a powerful method, it is difficult to keep track of the statistical errors involved in the individual histograms obtained in the simulations.

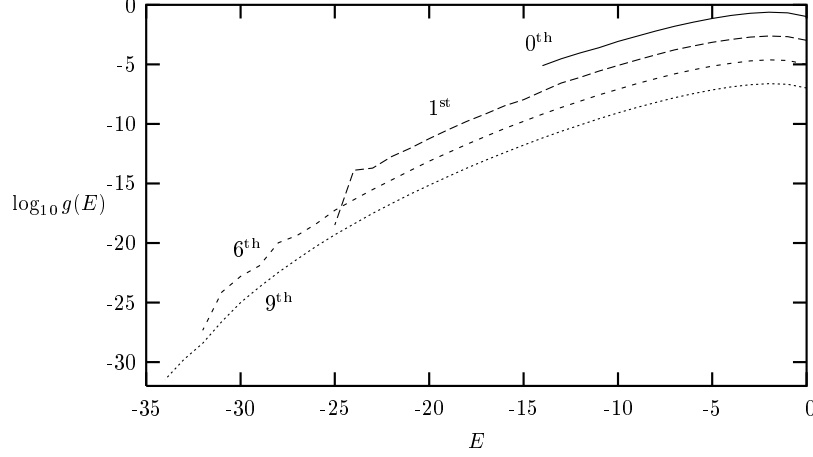
An alternative approach, in which the density of states  $g(E)$  is obtained within a *single* simulation without the necessity of a subsequent multi-histogram reweighting, is the combination of PERM with multicanonical sampling, the so-called multicanonical chain-growth method [34,35].

The general idea of multicanonical sampling [48,49] is to simulate the thermodynamic behavior of the system in a generalized (multicanonical) ensemble, where the energetic macrostates are distributed uniformly,  $p_{\text{muca}}(E) = \text{const}$ , which implies the introduction of multicanonical weight factors  $W_{\text{muca}}(E)$ . In typical multicanonical Monte Carlo simulations, the dynamics is therefore governed by a random walk in energy space. Hence, the sampling of entropically rare events is, in principle, as frequent as the sampling of highly degenerate energetic states. The acceptance probability for a new system configuration  $\mathbf{X}'$  with energy  $E'$  is  $w_{\text{muca}}(\mathbf{X} \rightarrow \mathbf{X}') = \min[1, \exp\{S(E(\mathbf{X}')) - S(E(\mathbf{X}))\}]$ , where  $S(E(\mathbf{X})) = -\ln W_{\text{muca}}(E(\mathbf{X}))$  is the microcanonical entropy. The canonical energy distribution  $p_{\text{can}}(E) \sim g(E) \exp(-E/k_B T)$  for a given temperature  $T$  is related with the multicanonical histogram via

$$p_{\text{can}}(E) \sim W_{\text{muca}}^{-1}(E) p_{\text{muca}}(E) e^{-E/k_B T}, \quad (8)$$

which implies that the multicanonical weights are proportional to the inverse density of states,  $W_{\text{muca}}(E) \sim g^{-1}(E)$ . Since  $g(E)$  is unknown, the determination of the weights  $W_{\text{muca}}(E)$  is not straightforward and must be performed in the first stage of the simulation in an iterative procedure [49].

The multicanonical extension of PERM requires two main changes compared to standard PERM. Firstly, the expression (6) for the weight factor is replaced by



**Fig. 5.** Estimates for the density of states  $g^{(i)}(E)$  for an exemplified heteropolymer with 42 monomers after several recursion levels. Since the curves would fall on top of each other, we have added, for better distinction, a suitable offset to the curves of the 1th, 6th, and 9th run. The estimate of the 0th run is normalized to unity.

$$W_n^{\text{MPERM}}(E_n) = W_{\text{muca},n}(E_n) \prod_{l=2}^n m_l, \quad W_{\text{muca},1} = 1, \quad (9)$$

where, according to multicanonical sampling, the multicanonical weight of the chain of current length  $n$  is related to the appropriate inverse density of states,  $W_{\text{muca},n}(E) \sim g_n^{-1}(E)$ . Note that the possibility to rewrite Eq. (9) in the recursive, factorized form

$$W_n^{\text{MPERM}}(E_n) = \prod_{l=2}^n m_l \frac{W_{\text{muca},l}(E_l)}{W_{\text{muca},l-1}(E_{l-1})} = W_{n-1}^{\text{MPERM}} m_n \frac{W_{\text{muca},n}(E_n)}{W_{\text{muca},n-1}(E_{n-1})} \quad (10)$$

is mainly responsible for the efficiency of this method as it ensures that rare-event (flat-histogram) sampling is performed in *all* intermediate steps of the growth process. This means that for a chain of length  $N$  all energy histograms are “flat”,  $H_n(E) \approx \text{const.}$  with  $n \leq N$ . The pruning-enriching scheme of PERM is completely carried over and remains unchanged with the exception that the thresholds (7) are re-expressed as

$$W_n^> = C_1 Z_n^{\text{MPERM}} \frac{c_n^2}{c_1^3}, \quad W_n^< = C_2 W_n^>, \quad (11)$$

i.e., in terms of the partition sum of the multicanonical ensemble,  $Z_n^{\text{MPERM}} = \sum_t W_n^{\text{MPERM}}(t)/c_1$ .

The second difference compared with the original PERM is the estimation of the multicanonical weights, as the densities of states  $g_n(E)$ ,  $n \leq N$ ,

are unknown in the beginning of the simulation. Therefore, the multicanonical weight factors  $W_{\text{muca},n}(E)$  must be determined iteratively for all stages  $n \leq N$  of the growth process [35]. The initial choice for the multicanonical weights is typically  $W_{\text{muca},n}^{(0)}(E) = 1 \forall n, E$ , making the zeroth recursion a pure PERM run at infinite temperature. The energy histograms are initialized with  $H_n^{(0)}(E) = 0$ . Performing the multicanonical chain growth according to the method described above, the histograms are accumulated by summing up the weights (10) of successively generated chains:

$$H_n^{(0)}(E) = \frac{1}{c_1} \sum_t W_{n,t}^{\text{MPERM}} \delta_{E_t E}, \quad (12)$$

where  $t$  labels the chain reaching length  $n$  in the growth process. Since this histogram is a first estimate for the density of states, the multicanonical weights for the following iteration are set to  $W_{\text{muca},n}^{(1)}(E) = 1/H_n^{(0)}(E)$ . Before starting the new recursion,  $Z_n^{\text{MPERM}}$ ,  $c_n$ ,  $W_n^<$  are reset to zero, and  $W_n^>$  to infinity (i.e. to the upper limit of the data type used to store this quantity). The iterative procedure is repeated until the weights  $W_{\text{muca},n}^{(i)}(E) = W_{\text{muca},n}^{(i-1)}(E)/H_n^{(i-1)}(E)$  are stabilized. In a long final production run  $i = I$ , the densities of states are then determined from

$$g_n^{(I)}(E) = \frac{H_n^{(I)}(E)}{W_{\text{muca},n}^{(I)}(E)}, \quad n \leq N. \quad (13)$$

For practical applications of this algorithm, in particular for studies of heteropolymers, it is more favorable to replace the original pruning-enrichment core, i.e., PERM [30], by the modern, improved variants nPERMss or nPERMis [33]. The combination of this more efficient chain-growth strategy with multicanonical sampling is straightforward. The details are explained in Refs. [34,35]. In Fig. 5, estimates for  $g^{(i)}(E)$  after the iterations  $i = 0, 1, 6$ , and 9 are shown for an exemplified heteropolymer with 42 monomers, whose thermodynamic properties will be discussed in more detail in Section 3.4. The zeroth run is a pure PERM estimate with a reliable precision over 5 orders of magnitude. As the simulations were effectively performed in the purely entropic regime at infinite temperature (i.e.,  $\beta = 0$ ), low-energy states are rarely sampled. In this case, chain growth is governed only by the weights (6) which are only products of free nearest-neighbor sites and therefore identical with the Rosenbluth weights for self-avoiding walks. The model-dependent energetic influence on the growth is thus irrelevant. The efficiency is improved in the successive recursions, where the multicanonical weights (10) are gradually refined and allow for a sampling of larger regions of the energy space. After only 10 recursions, the estimate for  $g(E)$  covers the whole accessible energy space (including the ground state) and ranges over 25 orders of magnitude.

It is worth noticing that the obtained density of states is *absolute*, i.e., estimates for the degeneracies of energetic states, in particular for the important

ground state, can directly be read off. Furthermore, the partition function is also absolutely estimated via  $Z = \sum_E g(E) \exp(-E/k_B T)$ . The reason is that these chain-growth methods perform a “biased” simple-sampling instead of importance sampling as is used in most Monte Carlo methods. With importance sampling, it is usually not possible to obtain an absolute estimate for  $g(E)$ .

The probability for energetic states in a canonical ensemble at temperature  $T$  is obtained from the density of states by Boltzmann reweighting via  $p_{\text{can}}(T) = g(E) \exp(-E/k_B T)/Z$ . Thus, statistical expectation values of energetic observables  $O(E)$  are simply given by  $\langle O \rangle = \sum_E O(E) p_{\text{can}}(E)$ . Thermal fluctuations of these quantities, e.g., defined by  $d\langle O \rangle/dT = (\langle O^2 \rangle - \langle O \rangle^2)/k_B T^2$ , are of particular interest for identifying temperature regions of thermodynamic activity. A very convenient measure for quantifying the cooperative behavior of a complex system is, e.g., the specific heat  $C_V = (\langle E^2 \rangle - \langle E \rangle^2)/k_B T^2$ .

### Decoupling energy scales: An instructive example

For systems, where different energy scales decouple, the density of states  $g(E)$  as a distribution of states with given *total* energy  $E$  is not the most useful quantity. As an important example, we consider the adsorption of a polymer to a substrate. In simple lattice models, only the number of intrinsic nearest-neighbor contacts between nonadjacent monomers,  $n_m$ , and the number of nearest-neighbor contacts of the polymer with the substrate,  $n_s$ , are counted. For the discussion of conformational transitions in the adsorption process later on, it is quite useful to rate intrinsic and binding forces against each other and therefore it is useful to introduce different energy scales  $\varepsilon_m$  and  $\varepsilon_s$  corresponding to the contact numbers  $n_m$  and  $n_s$ , respectively. A minimalistic model could then, for example, be defined by [50]

$$E(n_m, n_s) = -\varepsilon_m n_m - \varepsilon_s n_s \equiv -\varepsilon_0 (s n_m + n_s), \quad (14)$$

where the ratio  $s = \varepsilon_m/\varepsilon_s$  can be considered as kind of reciprocal solvent parameter (the larger  $s$ , the worse the quality of the solvent). The overall energy scale is simply  $\varepsilon_0 \equiv \varepsilon_s$ . Since the total energy  $E$  of the system depends on  $s$ , it would be necessary to fix its value in the previously described multicanonical chain-growth variant. Instead of determining the density of states  $g(E)$ , it would be more favorable to calculate the *contact density*  $g(n_m, n_s)$  which is independent of  $s$ . Knowing the contact density, the canonical probability for a system conformation with  $n_m$  monomer-monomer and  $n_s$  monomer-substrate contacts is given by  $p_{T,s}(n_m, n_s) = g(n_m, n_s) \exp[-E(n_m, n_s)/k_B T]/Z_{T,s}$ , where temperature  $T$  and solubility  $s$  are considered as fixed parameters. The statistical average of a quantity  $O(n_m, n_s)$  is then obtained as  $\langle O(n_m, n_s) \rangle = \sum_{n_m, n_s} O(n_m, n_s) p_{T,s}(n_m, n_s)$ . For the discussion of the conformational-phase diagram of the hybrid polymer-substrate system in solvent, it is useful to consider the dependence of fluctuations on temperature *and* solubility. As an

example, the specific heat can be expressed as

$$C_V(T, s) = k_B \left( \frac{\varepsilon_0}{k_B T} \right)^2 (s+1) \begin{pmatrix} \langle n_s^2 \rangle_c & \langle n_s n_m \rangle_c \\ \langle n_s n_m \rangle_c & \langle n_m^2 \rangle_c \end{pmatrix} \begin{pmatrix} s \\ 1 \end{pmatrix}, \quad (15)$$

where  $\langle xy \rangle_c = \langle xy \rangle - \langle x \rangle \langle y \rangle$  ( $x, y = n_m, n_s$ ) are the variances and covariances of the contact numbers. Note that the knowledge of  $g(n_m, n_s)$  enables reweighting of the specific heat to any pair of parameters  $T$  and  $s$ .

### Contact density chain-growth method

The determination of the contact density  $g(n_m, n_s)$  follows similar lines as the multicanonical chain-growth method for the estimation of the density of states. In fact, the only change in the algorithm described in the previous section is that the weights  $W_n^{\text{MPERM}}(E_n)$  defined in Eq. (10) are replaced by

$$W_n^{\text{CDPERM}}(n_m^{(n)}, n_s^{(n)}) = \prod_{l=2}^n m_l \frac{W_{\text{cd},l}(n_m^{(l)}, n_s^{(l)})}{W_{\text{cd},l-1}(n_m^{(l-1)}, n_s^{(l-1)})}, \quad (16)$$

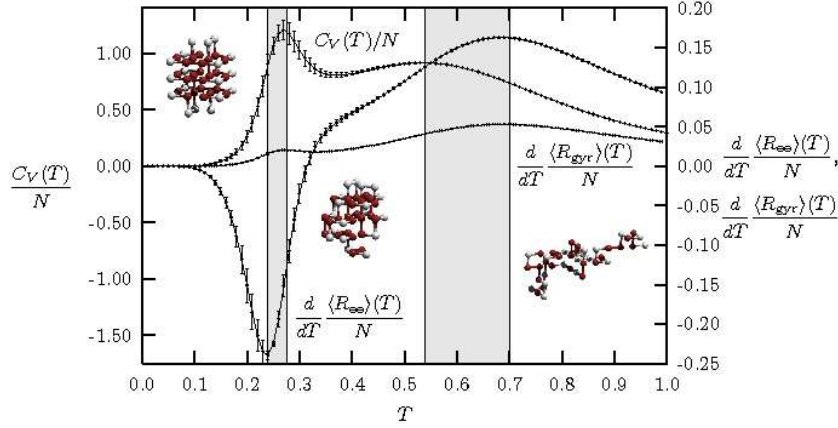
where the multi-contact weights  $W_{\text{cd},l}(n_m^{(l)}, n_s^{(l)}) \sim 1/g(n_m^{(l)}, n_s^{(l)})$  have again to be determined recursively.

The extension of this method incorporating more than two system parameters is straightforward, but the efficiency of flattening the high-dimensional histograms at all levels of the growth process decreases, whereas the storage requirements for these fields rapidly increase.

### 3.4 Bulk behavior of HP lattice proteins

Before embarking into the discussion of hybrid peptide-substrate systems, we investigate first the bulk behavior of HP peptides. In the folding process from a random-coil conformation towards the native fold, the protein experiences in many cases conformational transitions. These transitions typically require passing or circumventing of barriers in the free-energy landscape, which slows down the folding dynamics. Similar to thermodynamic phase transitions, conformational transitions can be identified by noticeable changes in the behavior of fluctuating quantities. Peaks and “shoulders” in specific-heat curves are, for example, typical signals for cooperative activity, because in the vicinity of the peak temperatures entropic changes separate qualitatively different classes of conformations (e.g., random coils and globular shapes). Since peptides are always of finite length due to their well-defined amino acid sequence, conformational transitions are not phase transitions in the strict thermodynamic sense. In consequence, fluctuations of different thermodynamic quantities typically do not exhibit the same peak structure, i.e., there is no “data collapse” which would allow the definition of a uniform transition temperature, where the phases are uniquely separated. Since for peptides different fluctuating





**Fig. 6.** Specific heat  $C_V$  and respective fluctuations of gyration radius and end-to-end distance,  $d\langle R_{\text{gyr}} \rangle/dT$  and  $d\langle R_{\text{ee}} \rangle/dT$ , as functions of temperature for the 42-mer.

quantities predict different transition temperatures, it is only possible to identify a temperature interval of conformational activity. This makes a precise quantitative analysis and a qualitative classification of such transitions difficult [34,35].

In the following, we discuss ground-state properties and thermodynamics for exemplified HP sequences. These results were obtained by employing the aforementioned multicanonical chain-growth method [34,35] for the standard version of the HP model [9]. An interesting example is the 42-monomer HP sequence representing the parallel  $\beta$ -helix protein *pectate lyase C* [51], which reads  $\text{PH}_2\text{PHPH}_2\text{PHPH}_2\text{H}_3\text{PHPH}_2\text{PHPH}_3\text{P}_2\text{HHPH}_2\text{PHPH}_2\text{P}$  [16]. Although it is not believed that specific protein properties such as the folding behavior and thermodynamics are conserved in a one-to-one transcription of an amino acid sequence into the hydrophobic-polar two-letter code, this example shows surprising coincidences of the real protein and the model, as the (low-degenerate) ground-state conformations in the HP model also exhibit two parallel helical segments. More interesting is, however, that the ground-state degeneracy is only  $g_0 = 4$  without trivial rotational symmetries [16]. With multicanonical chain-growth simulations [34], the ground-state degeneracy was precisely estimated as  $g_0 = 3.9 \pm 0.4$  [35]. In this simulation, 10 recursions were performed and in the production run, about  $5 \times 10^7$  chains entered into statistics.

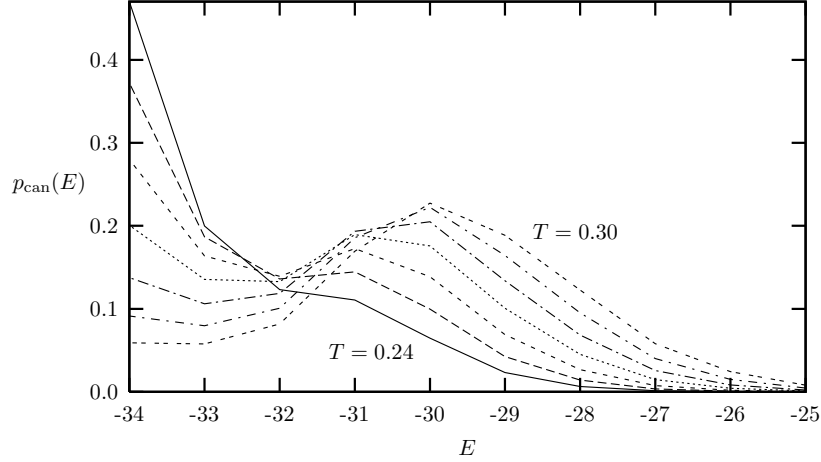
The low ground-state degeneracy is indeed remarkable, as it is extremely difficult to find *designing* sequences (which possess a nondegenerate ground state) with the standard HP model (cf. Table 1), in particular for comparatively long sequences. For a statistical analysis of the folding behavior of this 42-mer, the density of states, as has already been shown in Fig. 5, as well

as thermodynamic quantities and their fluctuations were calculated [34]. In Fig. 6, the specific heat  $C_V$  and fluctuations of the structural quantities radius of gyration  $R_{\text{gyr}}$  and end-to-end distance  $R_{\text{ee}}$  are plotted as functions of temperature. Two temperature regions of conformational activity (shaded in gray), where the curves of the fluctuating quantities exhibit extremal points, can clearly be separated.

For high temperatures, random conformations are favored. In consequence, in the corresponding, rather entropy-dominated ensemble, the high-degenerate high-energy structures govern the thermodynamic behavior of the macrostates. A typical representative is shown as an inset in the high-temperature pseudophase in Fig. 6. Annealing the system (or, equivalently, decreasing the solvent quality), the heteropolymer experiences a conformational transition towards globular macrostates. A characteristic feature of these intermediary “molten” globules is the compactness of the dominating conformations as expressed by a small gyration radius. Nonetheless, the conformations do not exhibit a noticeable internal long-range symmetry and behave rather like a fluid. Local conformational changes are not hindered by strong free-energy barriers. The situation changes by entering the low-temperature (or poor-solvent) conformational phase. In this region, energy dominates over entropy and the effectively attractive hydrophobic force favors the formation of a maximally compact core of hydrophobic monomers. Polar residues are expelled to the surface of the globule and form a shell that screens the core from the (fictitious) aqueous environment. In Fig. 7, we have plotted canonical energy distributions  $p_{\text{can}}(E)$  for several temperatures near the hydrophobic-core collapse transition. For temperatures above the transition region (which is between  $T^{(1)} = 0.24$  and  $T^{(2)} = 0.28$ , cf. Fig. 6), globular conformations are more probable, whereas for smaller temperatures hydrophobic-core states dominate. From the two-peak structure of the distributions in the transition region it can be concluded that this transition is first-order-like, i.e., both types of macrostates coexist in this temperature region.

The existence of the hydrophobic-core collapse renders the folding behavior of a heteropolymer different from crystallization or amorphous transitions of homopolymers. The reason is the disorder induced by the sequence of different monomer types. The hydrophobic-core formation is the main cooperative conformational transition which accompanies the tertiary folding process of a single-domain protein.

A very important aspect in the discussion of ground-state properties and conformational transitions towards the native fold is the influence of the heteropolymer sequence. For this purpose, we analyze ten designed sequences with 48 monomers, listed in Table 3, as given in Ref. [16]. The ratio between the numbers of hydrophobic and polar residues is one half for these HP proteins, i.e., the hydrophobicity is  $n_H = 24$ . The minimum energies we found from multicanonical chain-growth simulations [35] coincide with the values given in Refs. [16,33]. Also listed in Table 3 are the estimates for the degeneracies  $g_0$  of the respective ground-state energies. For comparison, previously



**Fig. 7.** Canonical energy distributions of the 42-mer for temperatures  $T = 0.24, 0.25, \dots, 0.30$  close to the hydrophobic-core collapse transition.

given lower bounds  $g_{\text{CHCC}}^<$  [37] are listed, which were obtained by means of the constraint-based hydrophobic core construction (CHCC) method [16]. Utilizing the idea of a highly compact hydrophobic core in the native fold, the hydrophobic monomers are in this method arranged in frames of maximal compactness. The number of the associated so-called Hamiltonian walks that connect the monomers, respecting the nonchangeable HP sequence, gives a lower bound for the degeneracy of the ground state. If, due to the sequence, a matching walk cannot be constructed, the compactness of the frame is relaxed and the search starts anew. The exact ground-state degeneracy would be obtained by scanning all frames and searching for conformations with the ground-state energy. Since the method is of exact enumeration type, the efforts of determining the precise ground-state degeneracy are enormous and, therefore, the main power of this method lies in the identification of native folds and the possibility to give a lower bound for its degeneracy.

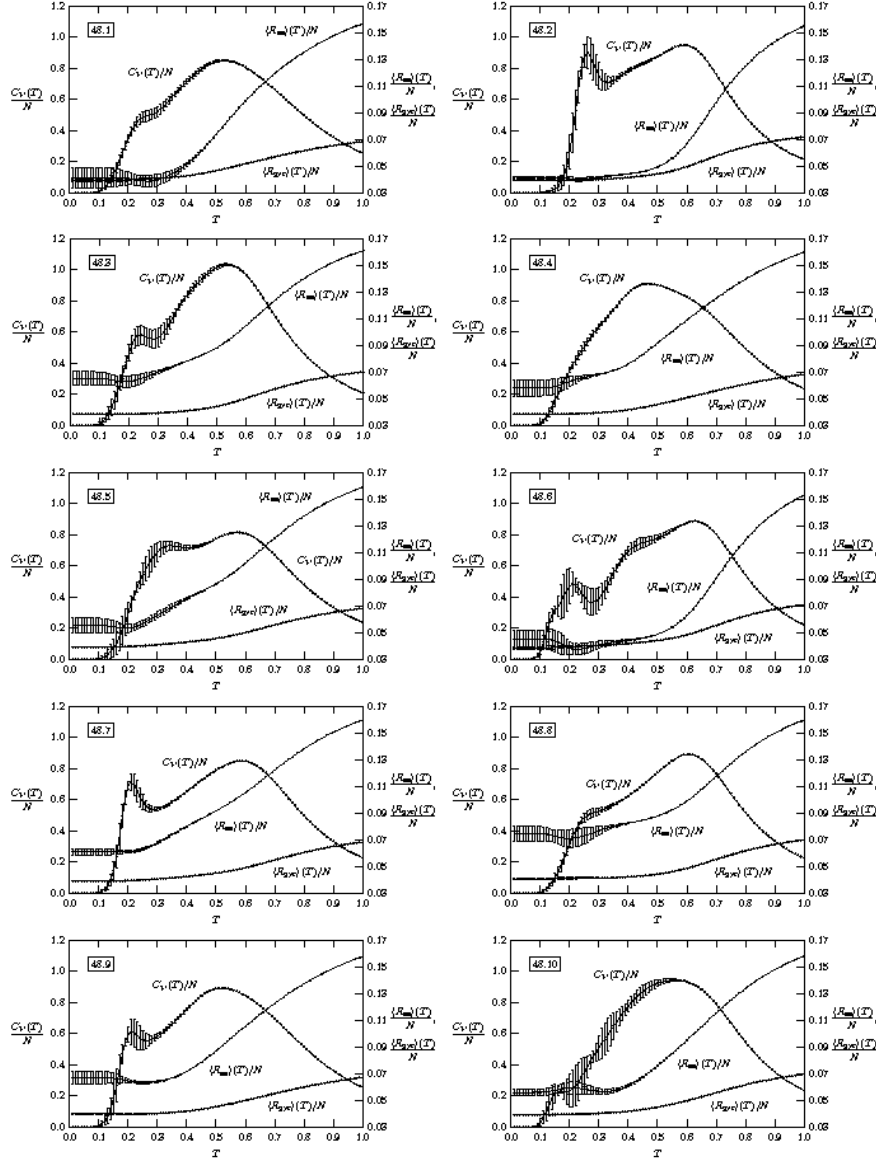
For the 48-mers, the  $g_0$  values obtained within the multicanonical chain-growth simulation lie indeed above these lower bounds or include it within the range of statistical errors. Notice that for the sequences 48.1, 48.5, and 48.8, the estimates for the ground-state degeneracy are much higher than the bounds  $g_{\text{CHCC}}^<$ . In these cases the smallest frame containing the entire hydrophobic core is rather large (cube containing  $4 \times 3 \times 3 = 36$  monomers with surface area  $A = 32$  [bond length]<sup>2</sup>) such that enumeration of this frame is cumbersome. For 48.5 and 48.8, we further found ground-state conformations lying in less compact frames (48.5:  $A = 32, 40, 42, 48, 52, 54$  [bond length]<sup>2</sup>, 48.8:  $A = 32, 40, 42$  [bond length]<sup>2</sup>) and those conformations would require still more effort to be identified with the CHCC algorithm, which was designed

**Table 3.** Ground-state energies  $E_{\min}$  and degeneracies  $g_0$  as estimated with the multicanonical chain-growth method [34,35] for ten HP sequences with 48 monomers. For comparison, we have also quoted the lower bounds on native degeneracies  $g_{\text{CHCC}}^<$  obtained by means of the CHCC (constraint-based hydrophobic core construction) method [16] as given in Ref. [37]. In both cases the constant factor 48 from rotational and reflection symmetries of conformations spreading into all three spatial directions was divided out.

No.	sequence	$E_{\min}$	$g_0 (\times 10^3)$	$g_{\text{CHCC}}^< (\times 10^3)$
48.1	HPH <sub>2</sub> P <sub>2</sub> H <sub>4</sub> PH <sub>3</sub> P <sub>2</sub> H <sub>2</sub> P <sub>2</sub> HPH <sub>3</sub> PHPH <sub>2</sub> P <sub>2</sub> H <sub>2</sub> P <sub>3</sub> HP <sub>8</sub> H <sub>2</sub>	-32	5226 ± 812	1500
48.2	H <sub>4</sub> PH <sub>2</sub> PH <sub>5</sub> P <sub>2</sub> HP <sub>2</sub> H <sub>2</sub> P <sub>2</sub> HP <sub>6</sub> HP <sub>2</sub> HP <sub>3</sub> HP <sub>2</sub> H <sub>2</sub> P <sub>2</sub> H <sub>3</sub> PH	-34	17 ± 8	14
48.3	PHPH <sub>2</sub> PH <sub>6</sub> P <sub>2</sub> HPHP <sub>2</sub> HPH <sub>2</sub> PHPH <sub>3</sub> HP <sub>2</sub> H <sub>2</sub> P <sub>2</sub> H <sub>2</sub> P <sub>2</sub> HPHP <sub>2</sub> HP	-34	6.6 ± 2.8	5.0
48.4	PHPH <sub>2</sub> P <sub>2</sub> HPH <sub>3</sub> P <sub>2</sub> H <sub>2</sub> PH <sub>2</sub> P <sub>3</sub> H <sub>5</sub> P <sub>2</sub> HPH <sub>2</sub> PHPH <sub>4</sub> HP <sub>2</sub> HPHP	-33	60 ± 13	62
48.5	P <sub>2</sub> HP <sub>3</sub> HPH <sub>4</sub> P <sub>2</sub> H <sub>4</sub> PH <sub>2</sub> PH <sub>3</sub> P <sub>2</sub> HPHPHP <sub>2</sub> HP <sub>6</sub> H <sub>2</sub> PH <sub>2</sub> PH	-32	1200 ± 332	54
48.6	H <sub>3</sub> P <sub>3</sub> H <sub>2</sub> PHPH <sub>2</sub> PH <sub>2</sub> PH <sub>2</sub> PHP <sub>7</sub> HPHP <sub>2</sub> HP <sub>3</sub> HP <sub>2</sub> H <sub>6</sub> PH	-32	96 ± 19	52
48.7	PHP <sub>4</sub> HPH <sub>3</sub> PHPH <sub>4</sub> PH <sub>2</sub> PH <sub>2</sub> P <sub>3</sub> HPHP <sub>3</sub> H <sub>3</sub> P <sub>2</sub> H <sub>2</sub> P <sub>2</sub> H <sub>2</sub> P <sub>3</sub> H	-32	58 ± 21	59
48.8	PH <sub>2</sub> PH <sub>3</sub> PH <sub>4</sub> P <sub>2</sub> H <sub>3</sub> P <sub>6</sub> HPH <sub>2</sub> P <sub>2</sub> H <sub>2</sub> PHP <sub>3</sub> H <sub>2</sub> PHPHPH <sub>2</sub> P <sub>3</sub>	-31	22201 ± 6594	306
48.9	PHPH <sub>4</sub> HPHPHP <sub>2</sub> HPH <sub>6</sub> P <sub>2</sub> H <sub>3</sub> PHP <sub>2</sub> HPH <sub>2</sub> P <sub>2</sub> HPH <sub>3</sub> P <sub>4</sub> H	-34	1.4 ± 0.5	1.0
48.10	PH <sub>2</sub> P <sub>6</sub> H <sub>2</sub> P <sub>3</sub> H <sub>3</sub> PHP <sub>2</sub> HPH <sub>2</sub> P <sub>2</sub> HP <sub>2</sub> HP <sub>2</sub> H <sub>2</sub> P <sub>2</sub> H <sub>7</sub> P <sub>2</sub> H <sub>2</sub>	-33	187 ± 87	188

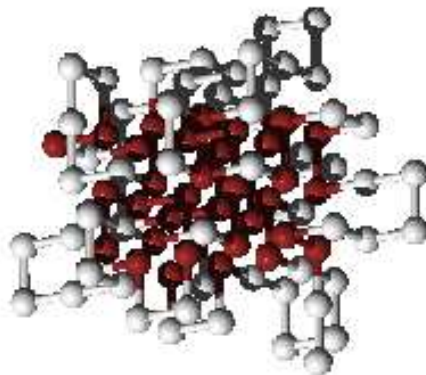
to locate global energy minima and therefore starts the search with the most compact hydrophobic frames. The ground-state energies of these examples are rather high ( $E_{\min} = -31$  for 48.8, and  $E_{\min} = -32$  for 48.1 and 48.5) and therefore a higher degeneracy seems to be natural. This is, however, only true, if there does not exist a conformational barrier that separates the compact H-core low-energy states from the general compact globules. Comparing the ground-state degeneracies and the low-temperature behavior of the specific heats for the sequences 48.1, 48.5, 48.6, and 48.7 (all of them having global energy minima with  $E_{\min} = -32$ ) as shown in Fig. 8, we observe that 48.6 and 48.7 with rather low ground-state degeneracy actually possess a pronounced low-temperature peak in the specific heat, while the higher-degenerate proteins 48.1 and 48.5 only show up a weak indication of a structural transition at low temperatures. The HP proteins 48.2, 48.3, and 48.9, which have the lowest minimum energy  $E_{\min} = -34$  among the examples in Table 3, have also the lowest ground-state degeneracies. These three candidates seem indeed to exhibit a rather strong ground-state – globule transition, as can be read off from the associated specific heats in Fig. 8.

In Fig. 8, also the mean end-to-end distances and mean radii of gyration are plotted as functions of temperature. Both quantities usually serve to interpret the conformational compactness of polymers. For HP proteins, the end-to-end distance is strongly influenced, however, by the types of monomers attached to the ends of the chain. It is easily seen from the figures that the 48mers with sequences starting and ending with a hydrophobic residue (48.1, 48.2, and 48.6) have a smaller mean end-to-end distance at low temperatures than the other examples from Table 3. The reason is that the ends can form hydrophobic contacts and therefore a reduction of the energy can be achieved. Thus, in these cases contacts between ends are usually favorable and the mean end-to-end distance is close to the mean radius of gyration.



**Fig. 8.** Specific heat, mean radius of gyration, and mean end-to-end distance for the ten 48-mers listed in Table 3 [35].

Interestingly, there exists indeed a crossover region, where  $\langle R_{ee} \rangle < \langle R_{\text{gyr}} \rangle$ . Comparing with the behavior of the specific heat, this interval is close to the region, where the phase dominated by low-energy states crosses over to the globule-favored phase. The hydrophobic contact between the ends is strong



**Fig. 9.** Compact hydrophobic-core conformation of the 103-mer [35] used in the peptide adsorption study [68,69]. Dark spheres correspond to hydrophobic monomers and light spheres mark polar residues.

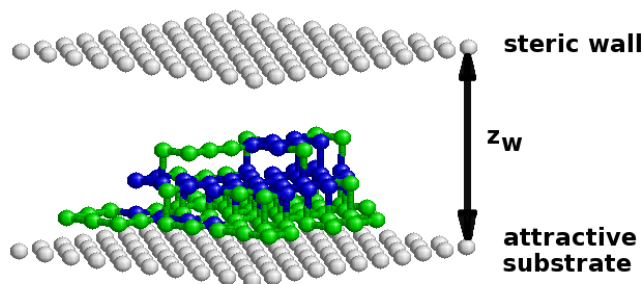
enough to resist the thermal fluctuations in that temperature interval. The reason is that, once such a hydrophobic contact between the ends is established, usually other in-chain hydrophobic monomers are attracted and form a hydrophobic core surrounding the end-to-end contact. Thus, before the contact between the ends is broken, an increase of the temperature first leads to a melting of the surrounding contacts. The entropic freedom to form new conformations is large since the low-energy states are all relatively high degenerate and do not possess symmetries requiring an appropriate amount of heat to be broken. For sequences possessing mixed or purely polar ends, the mean end-to-end distance and mean radius of gyration differ much stronger, as there is no energetic reason, why the ends occupy nearest-neighbor positions.

In conclusion, we see that for longer chains the strength of the low-temperature transition not only depends on low ground-state degeneracies as it does for short chains [34]. Rather, the influence of the higher-excited states cannot be neglected. A striking example is sequence 48.4 with rather low ground-state degeneracy, but only weak signals for a low-temperature transition.

### 3.5 Specificity of protein adsorption to selective solid substrates

In this section, we discuss results of a simple lattice model similar to Eq. (14) for analyzing the conformational behavior of HP proteins in adsorption processes to different, specific solid substrates. The objective is the determination of a pseudophase diagram, which allows for the classification of conformational subphases in dependence of the external parameters temperature and solubility of the surrounding (implicit) solvent.

The recent developments in single molecule experiments at the nanometer scale, e.g., by means of atomic force microscopy (AFM) [52] and optical



**Fig. 10.** Lattice model used in the peptide-substrate adsorption study.

tweezers [53], allow now for a more detailed exploration of structural properties of polymers in the vicinity of adsorbing substrates [54]. The possibility to perform such studies is of essential biological and technological significance. From the biological point of view the understanding of the binding and docking mechanisms of proteins at cell membranes is important for the reconstruction of biological cell processes. Similarly, specificity of peptides and binding affinity to selected substrates could be of great importance for future electronic nanoscale circuits and pattern recognition nanosensory devices [55]. The study of hybrid interface models has considerable applications for a broad variety of problems, e.g., understanding the mechanisms of protein–ligand binding [56], prewetting and layering transitions in polymer solutions as well as dewetting of polymer films [57,58], molecular pattern recognition [59], electrophoretic polymer deposition and growth [60]. Recently, the influence of adhesion and steric hindrance for polymers grafted to a flat substrate [61,50,62–65], conformational pseudophase transitions for nongrafted polymers and peptides in a cavity with attractive substrate [66–69], the shape response to pulling forces [70,71] or external fields [72] were subject of computer simulations and analytical approaches of different models. The question how a flexible substrate, e.g., a cell membrane, bends as a reaction of a grafted polymer, was, for example, addressed in Ref. [73]. Proteins exhibit a strong specificity as the affinity of peptides to adsorb at surfaces depends on the amino acid sequence, solvent properties, and substrate shape. This was experimentally and numerically studied, e.g., for peptide-metal [74,75] and peptide-semiconductor [76,77] interfaces. Binding/folding and docking properties of lattice heteropolymers at an adsorbing surface were also subject of numerical studies [78].

### Lattice model for hybrid peptide-substrate interfaces

For the study of hybrid peptide-substrate models, we use the HP transcription of the 103-residue protein *cytochrome c*, which was extensively studied in the past [33,35,38,39]. The HP sequence contains 37 hydrophobic and 66 polar residues. A conformation with a highly compact hydrophobic-core, ex-

hibiting 56 hydrophobic contacts, is shown in Fig. 9. This lattice peptide resides in a cavity with an attractive substrate (see Fig. 10). For regularization of the upper halfspace, an additional steric wall in a distance  $z_w$  is introduced. The value of  $z_w$  is chosen sufficiently large to keep the influence on the unbound heteropolymer small (in the actual example,  $z_w = 200$  was used). In order to study the specificity of residue binding, we distinguish three substrates with different affinities to attract the peptide monomers: (a) the type-independent attractive, (b) the hydrophobic, and (c) the polar substrate. The number of corresponding nearest-neighbor contacts between monomers and substrate shall be denoted as  $n_s^{H+P}$ ,  $n_s^H$ , and  $n_s^P$ , respectively. In analogy to the polymer-substrate model (14), we express the energy of the hybrid peptide-substrate system simply by

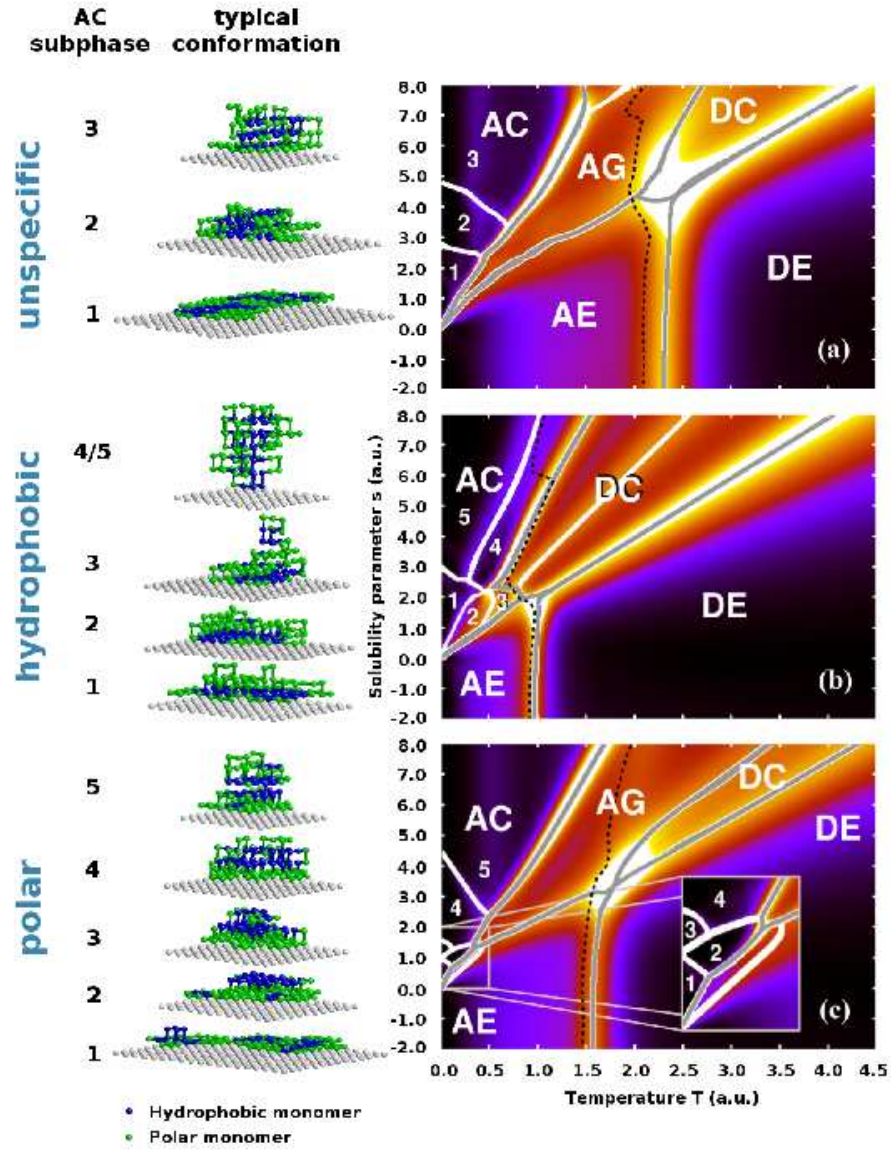
$$E_s(n_s, n_{\text{HH}}) = -\varepsilon_0(n_s + sn_{\text{HH}}), \quad (17)$$

where  $n_s = n_s^{H+P}$ ,  $n_s^P$ , or  $n_s^H$  depending on the substrate (we set  $\varepsilon_0 = 1$  in the following). The solubility (or reciprocal solvent parameter)  $s$  is, as well as the temperature  $T$ , an external parameter. It controls the quality of the solvent (the larger the value of  $s$ , the worse the solvent). This model was investigated by means of the contact-density chain-growth algorithm (see Sec. 3.3), which allows a direct estimation of the degeneracy (or density)  $g(n_s, n_{\text{HH}})$  of macrostates of the system with given contact numbers  $n_s$  and  $n_{\text{HH}}$  [68,69]. In contrast to move-set based Metropolis Monte Carlo or conventional chain-growth methods which would require many separate simulations to obtain results for different parameter pairs  $(T, s)$  and which frequently suffer from slowing down in the low-temperature sector, the contact-density chain-growth method allows the computation of the *complete* contact density for each system within a *single* simulation run. Since the contact density is independent of temperature and solubility, energetic quantities such as the specific heat (15) can easily be calculated for all values of  $T$  and  $s$ . Nonenergetic quantities require accumulated densities to be measured within the simulation, but this is also no problem.

### Conformational adsorption behavior in dependence of temperature and solubility

In Figs. 11(a)–(c) the color-coded profiles of the specific heats for the different substrates are shown (the brighter the larger the value of  $C_V$ ). We interpret the ridges (for accentuation marked by white and gray lines) as the boundaries of the pseudophases. The gray lines indicate the main transition lines, while the white lines separate pseudophases that strongly depend on specific properties of the heteropolymer, such as its exact number and sequence of hydrophobic and polar monomers. With its degeneracy  $g(n_s, n_{\text{HH}})$ , we define the contact free energy as  $F_{T,s}(n_s, n_{\text{HH}}) = E_s(n_s, n_{\text{HH}}) - T \ln g(n_s, n_{\text{HH}})$  and the probability for a macrostate with  $n_s$  substrate and  $n_{\text{HH}}$  hydrophobic contacts

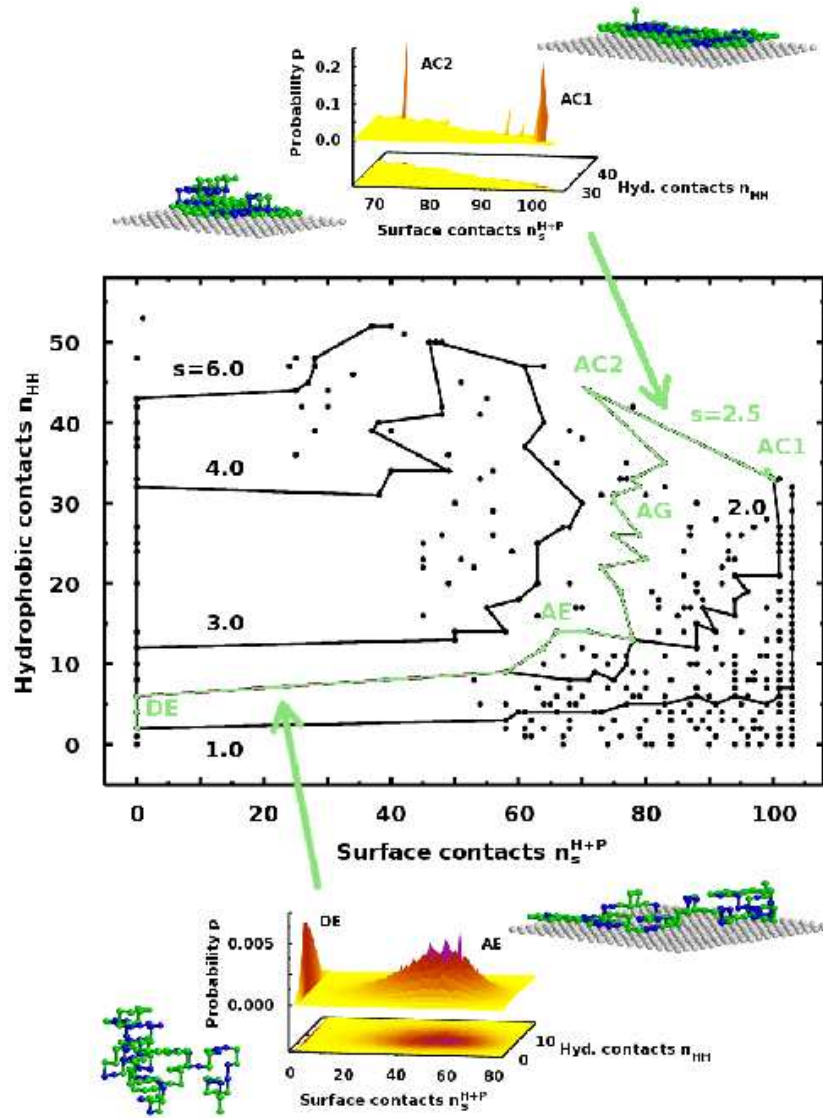




**Fig. 11.** Specific-heat profiles as a function of temperature  $T$  and solubility parameter  $s$  of the 103-mer near three different substrates that are attractive for (a) all, (b) only hydrophobic, and (c) only polar monomers. White lines indicate the ridges of the profile. Gray lines mark the main "phase boundaries". The dashed black line represents the first-order-like binding/unbinding transition state, where the contact free energy possesses two minima (the adsorbed and the desorbed state). In the left panel typical conformations dominating the associated AC phases of the different systems are shown.

as  $p_{T,s}(n_s, n_{\text{HH}}) \sim g(n_s, n_{\text{HH}}) \exp(-E_s/T)$ . Assuming that the minimum of the free-energy landscape  $F_{T,s}(n_s^{(0)}, n_{\text{HH}}^{(0)}) \rightarrow \min$  for given external parameters  $s$  and  $T$  is related to the class of macrostates with  $n_s^{(0)}$  surface and  $n_{\text{HH}}^{(0)}$  hydrophobic contacts, this class dominates the phase the system resides in. For this reason, it is instructive to calculate all minima of the contact free energy and to determine the associated contact numbers in a wide range of values for the external parameters.

The map of all possible free-energy minima in the range of external parameters  $T \in [0, 10]$  and  $s \in [-2, 10]$  is shown in Fig. 12 for the peptide in the vicinity of a substrate that is equally attractive for both hydrophobic and polar monomers. Solid lines visualize “paths” through the free energy landscape when changing temperature under constant solvent ( $s = \text{const}$ ) conditions. Let us follow the exemplified trajectory for  $s = 2.5$ . Starting at very low temperatures, we know from the pseudophase diagram in Fig. 11(a) that the system resides in pseudophase AC1. This means that the macrostate of the peptide is dominated by the class of compact, film-like single-layer conformations. The system obviously prefers surface contacts at the expense of hydrophobic contacts. Nonetheless, the formation of compact hydrophobic domains in the two-dimensional topology is energetically favored but maximal compactness is hindered by the steric influence of the substrate-binding polar residues. Increasing the temperature, the system experiences close to  $T \approx 0.35$  a sharp first-order-like conformational transition, and a second layer forms (AC2). This is a mainly entropy-driven transition as the extension into the third dimension perpendicular to the substrate surface increases the number of possible peptide conformations. Furthermore, the loss of energetically favored substrate contacts of polar monomers is partly compensated by the energetic gain due to the more compact hydrophobic domains. Increasing the temperature further, the density of the hydrophobic domains reduces and overall compact conformations dominate in the globular pseudophase AG. Reaching AE, the number of hydrophobic contacts decreases further, and also the total number of substrate-contacts. Extended, dissolved conformations dominate. The transitions from AC2 to AE via AG are comparatively “smooth”, i.e., no immediate changes in the contact numbers passing the transition lines are noticed. Therefore, these conformational transitions could be classified as second-order-like. The situation is different when approaching the unbinding transition line from AE close to  $T \approx 2.14$ . This transition is accompanied by a dramatic loss of substrate contacts – the peptide desorbs from the substrate and behaves in pseudophase DE like a free peptide, i.e., the substrate and the opposite neutral wall regularize the translational degree of freedom perpendicular to the walls, but rotational symmetries are unbroken (at least for conformations not touching one of the walls). As the probability distribution in Fig. 12 shows, the unbinding transition is also first-order-like, i.e., close to the transition line, there is a coexistence of adsorbing and desorbing classes of conformations.

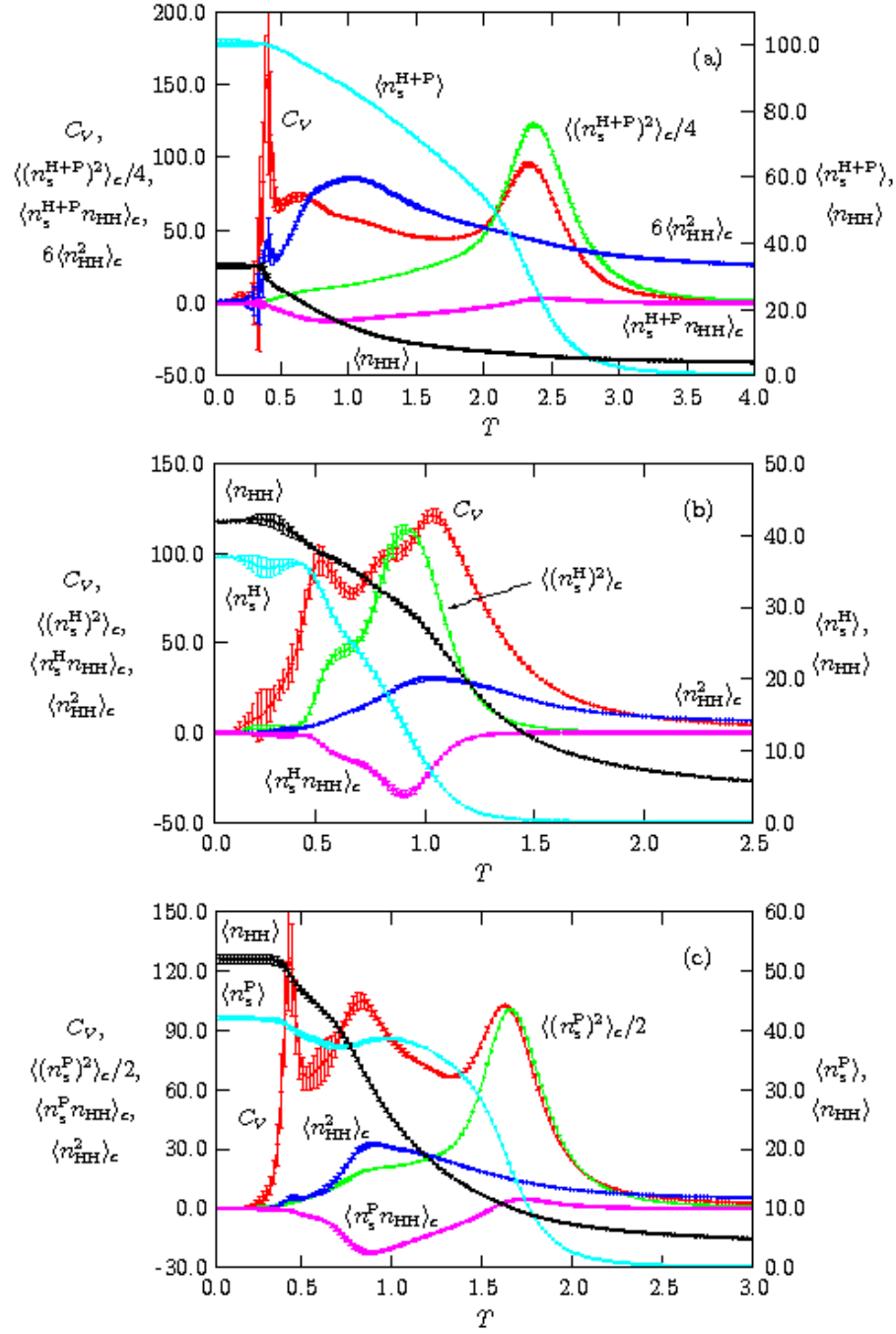


**Fig. 12.** Contact-number map of all free-energy minima for the 103-mer and substrate equally attractive to all monomers. Full circles correspond to minima of the contact free energy  $F_{T,s}(n_s^{H+P}, n_{HH})$  in the parameter space  $T \in [0, 10]$ ,  $s \in [-2, 10]$ . Lines illustrate how the contact free energy changes with the temperature at constant solvent parameter  $s$ . For the exemplified solvent with  $s = 2.5$ , the peptide experiences near  $T = 0.35$  a sharp first-order-like layering transition between single- to double-layer conformations (AC1,2). Passing the regimes of adsorbed globules (AG) and expanded conformations (AE), the discontinuous binding/unbinding transition from AE to DE happens near  $T = 2.14$ . In the DE phase the ensemble is dominated by desorbed, expanded conformations. Representative conformations of the phases are shown next to the respective peaks of the probability distributions.

Despite the surprisingly rich and complex phase behavior there are main “phases” that can be distinguished in all three systems. These are separated in Figs. 11(a)–(c) by gray lines. Comparing the three systems we find that they all possess pseudophases, where adsorbed compact (AC), adsorbed expanded (AE), desorbed compact (DC), and desorbed expanded (DE) conformations dominate. “Compact” here means that the heteropolymer has formed a dense hydrophobic core, while expanded conformations are dissolved, random-coil like. The sequence and substrate specificity of heteropolymers generates, of course, a rich set of new interesting and selective phenomena not available for homopolymers. One example is the pseudophase of adsorbed globules (AG), which is noticeably present only in those systems, where all monomers are equally attractive to the substrate (Fig. 11(a)) and where polar monomers favor contact with the surface (Fig. 11(b)). In this phase, the conformations are intermediates in the binding/unbinding region. This means that monomers currently desorbed from the substrate have not yet found their position within a compact conformation. Therefore, the hydrophobic core, which is smaller than in the respective adsorbed phase (i.e., at constant solubility  $s$ ), appears as a loose cluster of hydrophobic monomers.

In Figs. 13(a)–(c), we have plotted, exemplified for  $s = 2$ , the statistical averages of the contact numbers  $n_s$  and  $n_{HH}$  as well as their variances and covariances for the three systems. For comparison we have also included the specific heat, whose peaks correspond to the intersected transition lines of Figs. 11(a)–(c) at  $s = 2$ . From Figs. 13(a) and (c) we read off that the transition from AC to AG near  $T \approx 0.4$  is mediated by fluctuations of the intrinsic hydrophobic contacts. The very dense hydrophobic domains in the AC subphases lose their compactness. This transition is absent in the hydrophobic-substrate system (Fig. 13(b)). The signal seen belongs to a hydrophobic layering AC subphase transition, which influences mainly the number of surface contacts  $n_s^H$ . The second peak of the specific heats belongs to the transition between adsorbed compact or globular (AC/AG) and expanded (AE) conformations. This behavior is similar in all three systems. Remarkably, it is accompanied by a strong anti-correlation between surface and intrinsic contact numbers,  $n_s$  and  $n_{HH}$ . Not surprisingly, the hydrophobic contact number  $n_{HH}$  fluctuates stronger than the number of surface contacts, but apparently in a different way. Dense conformations with hydrophobic core (and therefore many hydrophobic contacts) possess a relatively small number of surface contacts. Vice versa, conformations with many surface contacts cannot form compact hydrophobic domains. Finally, the third specific heat peak marks the binding/unbinding transition, which is, as expected, due to a strong fluctuation of the surface contact number.

The strongest difference between the three systems is their behavior in pseudophase AC, which is roughly parameterized by  $s > 5T$ . If hydrophobic and polar monomers are equally attracted by the substrate (Fig. 11(a)), we find three AC subphases in the parameter space plotted. In subphase AC1, film-like conformations dominate, i.e., all 103 monomers are in contact with

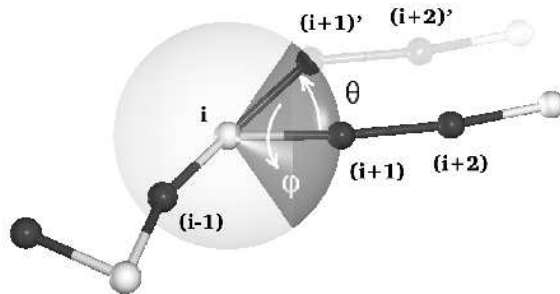


**Fig. 13.** Temperature dependence of specific heat, correlation matrix components, and contact number expectation values of the 103mer for surfaces attractive for (a) all, (b) only hydrophobic, and (c) only polar monomers at  $s = 2$ .

the substrate. Due to the good solvent quality in this region, the formation of a hydrophobic core is less attractive than the maximal deposition of all monomers at the surface, the ground state is  $(n_s^{\text{H+P}}, n_{\text{HH}})_{\min} = (103, 32)$ . In fact, instead of a single compact hydrophobic core there are nonconnected hydrophobic clusters. At least on the used simple cubic lattice and the chosen sequence, the formation of a single hydrophobic core is necessarily accompanied by an unbinding of certain polar monomers and, in consequence, an extension of the conformation into the third spatial dimension. In fact, this happens when entering AC2  $[(n_s^{\text{H+P}}, n_{\text{HH}})_{\min} = (64, 47)]$ , where a single hydrophobic two-layer domain has formed at the expense of losing surface contacts. In AC3, the heteropolymer has maximized the number of hydrophobic contacts and only local arrangements of monomers on the surface of the very compact structure lead to the still possible maximum number of substrate contacts.  $F_{T,s}$  is minimal for  $(n_s^{\text{H+P}}, n_{\text{HH}})_{\min} = (40, 52)$ .

The behavior of the heteropolymer adsorbed at a surface that is only attractive to hydrophobic monomers (Fig. 11(b)) is apparently different in the AC phase. Since surface contacts of polar monomers are energetically not favored, the subphase structure is determined by the competition of two hydrophobic forces: substrate attraction and formation of intrinsic contacts. In AC1, the number of hydrophobic substrate contacts is maximal for the single hydrophobic layer,  $(n_s^{\text{HH}}, n_{\text{HH}})_{\min} = (37, 42)$ . The *single* two-dimensional hydrophobic domain is also maximally compact, at the expense of displacing polar monomers into a second layer. In subphase AC2 intrinsic contacts are entropically broken with minimal free energy for  $35 \leq n_{\text{HH}} \leq 40$ , while  $n_s^{\text{HH}} = 37$  remains maximal. Another AC subphase, AC3, exhibits a hydrophobic layering transition at the expense of hydrophobic substrate contacts. Much more interesting is the subphase transition from AC1 to AC5. The number of hydrophobic substrate contacts  $n_s^{\text{HH}}$  of the ground-state conformation dramatically decreases (from 37 to 4) and the hydrophobic monomers collapse in a one-step process from the compact two-dimensional domain to the maximally compact three-dimensional hydrophobic core. The conformations are mushroom-like structures grafted at the substrate. AC4 is similar to AC5, with advancing desorption.

Not less exciting is the subphase structure of the heteropolymer interacting with a polar substrate (Fig. 11(c)). For small values of  $s$  and  $T$ , the behavior of the heteropolymer is dominated by the competition between polar monomers contacting the substrate and hydrophobic monomers favoring the formation of a hydrophobic core, which, however, also requires cooperativity of the polar monomers. In AC1, film-like conformations  $(n_s^{\text{P}}, n_{\text{HH}})_{\min} = (66, 31)$  with disconnected hydrophobic clusters dominate. Entering AC2, hydrophobic contacts are energetically favored and a second hydrophobic layer forms at the expense of a reduction of polar substrate contacts  $[(n_s^{\text{P}}, n_{\text{HH}})_{\min} = (61, 37)]$ . In AC3, the upper layer is mainly hydrophobic  $[(n_s^{\text{P}}, n_{\text{HH}})_{\min} = (53, 45)]$ , while the poor quality of the solvent ( $s$  large) and the comparatively strong hydrophobic force let the conformation further collapse [AC4:  $(n_s^{\text{P}}, n_{\text{HH}})_{\min} = (42, 52)]$



**Fig. 14.** Spherical update of the bond vector between the  $i$ th and  $(i+1)$ th monomer.

and the steric cooperativity forces more polar monomers to break the contact to the surface and to form a shell surrounding the hydrophobic core  $[(n_s^P, n_{HH})_{\min} = (33, 54)]$  in AC5].

## 4 Going off-lattice: Folding behavior of heteropolymers in the AB continuum model

The lattice models discussed in the previous sections suffer from the fact that the results for the finite-length heteropolymers typically depend on the underlying lattice type. It is difficult to separate realistic effects from artefacts induced by the use of a certain lattice structure. This problem can be avoided, in principle, by studying off-lattice heteropolymers, where the degrees of freedom are continuous. On the other hand, this advantage is partly counterbalanced by the increasing computational efforts for sampling the relevant regions of the conformational state space. In consequence, a precise analysis of statistical properties of off-lattice heteropolymers by means of generalized-ensemble methods can reliably be performed only for chains much shorter than those considered in the lattice studies. In the following, we focus on hydrophobic-polar heteropolymers with 20 monomers employing the so-called AB model [79], where  $A$  monomers are hydrophobic and residues of type  $B$  are polar (or hydrophilic).

### 4.1 Modeling and updating

We denote the spatial position of the  $i$ th monomer in a heteropolymer consisting of  $N$  residues by  $\mathbf{r}_i$ ,  $i = 1, \dots, N$ , and the vector connecting nonadjacent monomers  $i$  and  $j$  by  $\mathbf{r}_{ij}$ . For covalent bond vectors, we set  $|\mathbf{b}_i| \equiv |\mathbf{r}_{i,i+1}| = 1$ . The bending angle between monomers  $k$ ,  $k+1$ , and  $k+2$  is  $\vartheta_k$  ( $0 \leq \vartheta_k \leq \pi$ ) and  $\sigma_i = A, B$  symbolizes the type of the monomer. In the AB model [79], the energy of a conformation is given by

$$E = \frac{1}{4} \sum_{k=1}^{N-2} (1 - \cos \vartheta_k) + 4 \sum_{i=1}^{N-2} \sum_{j=i+2}^N \left( \frac{1}{r_{ij}^{12}} - \frac{C(\sigma_i, \sigma_j)}{r_{ij}^6} \right), \quad (18)$$

where the first term is the bending energy and the sum runs over the  $(N-2)$  bending angles of successive bond vectors. The second term partially competes with the bending barrier by a potential of Lennard-Jones type. It depends on the distance between monomers being non-adjacent along the chain and accounts for the influence of the AB sequence on the energy. The long-range behavior is attractive for pairs of like monomers and repulsive for  $AB$  pairs of monomers:

$$C(\sigma_i, \sigma_j) = \begin{cases} +1, & \sigma_i, \sigma_j = A, \\ +1/2, & \sigma_i, \sigma_j = B, \\ -1/2, & \sigma_i \neq \sigma_j. \end{cases} \quad (19)$$

The Monte Carlo simulation of this model is not straightforward as strictly local updates are not possible. A simple nonlocal update of a given conformation can be performed by using the procedure displayed in Fig. 14. Since the length of the bonds is fixed ( $|\mathbf{b}_k| = 1, k = 1, \dots, N-1$ ), the  $(i+1)$ th monomer lies on the surface of a sphere with radius unity around the  $i$ th monomer. Therefore, spherical coordinates are the natural choice for calculating the new position of the  $(i+1)$ th monomer on this sphere. For the reason of efficiency, we do not select any point on the sphere but restrict the choice to a spherical cap with maximum opening angle  $2\theta_{\max}$  (the dark area in Fig. 14). Thus, to change the position of the  $(i+1)$ th monomer to  $(i+1)'$ , we select the angles  $\theta$  and  $\varphi$  randomly from the respective intervals  $\cos \theta_{\max} \leq \cos \theta \leq 1$  and  $0 \leq \varphi \leq 2\pi$ , which ensure a uniform distribution of the  $(i+1)$ th monomer positions on the associated spherical cap. After updating the position of the  $(i+1)$ th monomer, the following monomers in the chain are simply translated according to the corresponding bond vectors which remain unchanged in this type of update. Only the bond vector between the  $i$ th and the  $(i+1)$ th monomers is rotated, all others keep their direction. This is similar to single spin updates in local-update Monte Carlo simulations of the classical Heisenberg model with the difference that in addition to local energy changes long-range interactions of the monomers, changing their relative position to each other, have to be computed anew after the update. For simulations in the state space of dense conformations it is recommendable to choose a rather small opening angle, e.g.,  $\cos \theta_{\max} = 0.99$ , in order to be able to sample also very narrow and deep valleys in the landscape of angles.

For the following discussion of folding channels of 20mers [85], these updates were used in combination with multicanonical sampling [48,49].

## 4.2 Characteristic protein folding channels and free-energy landscapes from coarse-grained modeling

The folding process of proteins is necessarily accompanied by cooperative conformational changes. Although not phase transitions in the strict sense, it



**Table 4.** The three AB 20-mers studied and the values of the associated (putative) global energy minima. Note that the given values for sequence S3 belong to two different, almost degenerate folds (cf. Fig. 16).

label	sequence	global energy minimum [85]
S1	$BA_6BA_4BA_2BA_2B_2$	-33.8236
S2	$A_4BA_2BABA_2B_2A_3BA_2$	-34.4892
S3	$A_4B_2A_4BA_2BA_3B_2A$	-33.5838, -33.5116

should be expected that one or a few parameters can be defined that enable the description of the structural ordering process [80,81]. The number of degrees of freedom in most all-atom models is given by the dihedral torsional backbone and side-chain angles. In coarse-grained  $C^\alpha$  models as the AB model used in this study, the original dihedral angles are replaced by a set of virtual torsional and bond angles. In fact, the number of degrees of freedom is not necessarily reduced in simplified off-lattice models. Therefore, the complexity of the space of degrees of freedom is comparable with more realistic models, and it is also a challenge to identify a suitable order parameter for the folding in such minimalistic heteropolymer models.

In analogy to studies of the specific folding behavior in all-atom protein models [82,83], it is suitable to define a generalized variant of an angular overlap order parameter as introduced in Ref. [84]. The idea is to define a simple and computationally low-cost measure for the similarity of two conformations, where the differences of the angular degrees of freedom are calculated. In order to consider this parameter as kind of order parameter, it is useful to compare conformations  $\mathbf{X} = (\mathbf{r}_1, \dots, \mathbf{r}_N)$  of the actual ensemble with a suitable reference conformation, which is preferably chosen to be the global-energy minimum conformation  $\mathbf{X}^{(0)}$ . We define the overlap parameter as follows:

$$Q(\mathbf{X}) = 1 - d(\mathbf{X}). \quad (20)$$

With  $N_b = N - 2$  and  $N_t = N - 3$  being the respective numbers of bond angles  $\Theta_i$  and torsional angles  $\Phi_i$ , the angular deviation between the conformations is calculated according to

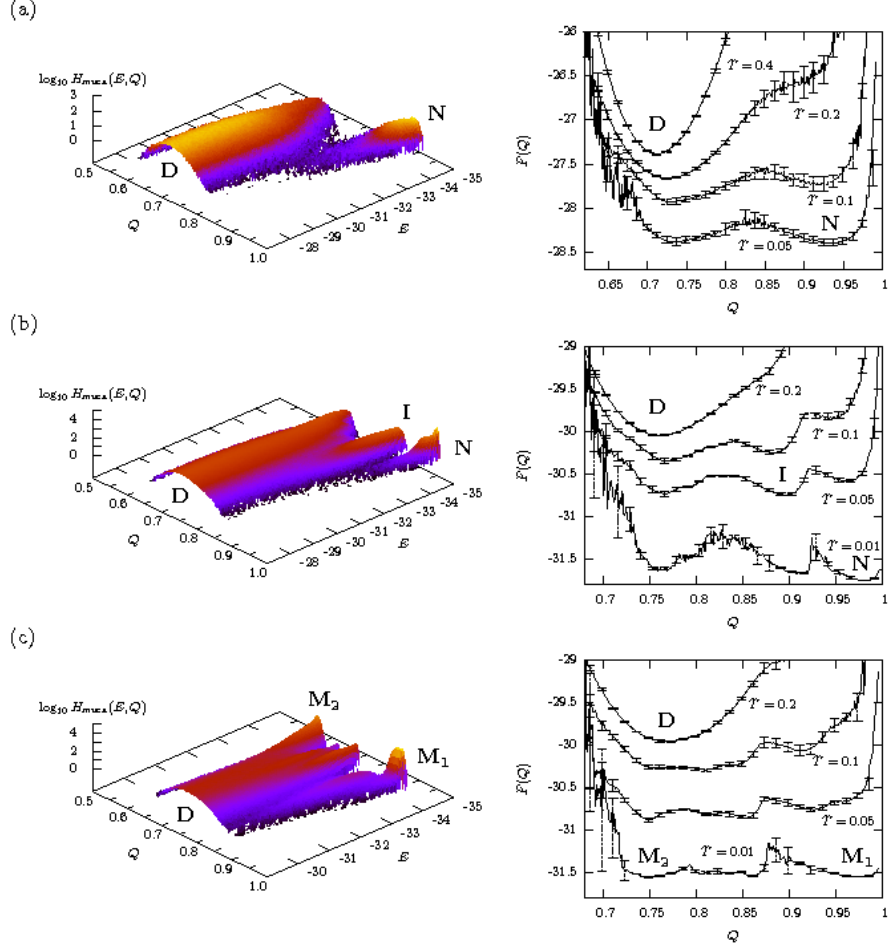
$$d(\mathbf{X}) = \frac{1}{\pi(N_b + N_t)} \left[ \sum_{i=1}^{N_b} d_b(\Theta_i) + \min_{r=\pm} \left( \sum_{i=1}^{N_t} d_t^r(\Phi_i) \right) \right], \quad (21)$$

where

$$d_b(\Theta_i) = |\Theta_i - \Theta_i^{(0)}|, \quad (22)$$

$$d_t^\pm(\Phi_i) = \min \left( |\Phi_i \pm \Phi_i^{(0)}|, 2\pi - |\Phi_i \pm \Phi_i^{(0)}| \right). \quad (23)$$

Here it is taken into account that the AB model is invariant under the reflection symmetry  $\Phi_i \rightarrow -\Phi_i$ . Thus, it is not useful to distinguish between



**Fig. 15.** Multicanonical histograms  $H_{\text{muca}}(E, Q)$  of energy  $E$  and angular overlap parameter  $Q$  and free-energy landscapes  $F(Q)$  at different temperatures for the three sequences (a) S1, (b) S2, and (c) S3. The reference folds reside at  $Q = 1$  and  $E = E_{\text{min}}$  [85].

reflection-symmetric conformations and therefore only the larger overlap is considered. Since  $-\pi \leq \Phi_i \leq \pi$  and  $0 \leq \Theta_i \leq \pi$ , the overlap is unity, if all angles of the conformations  $\mathbf{X}$  and  $\mathbf{X}^{(0)}$  coincide, else  $0 \leq Q < 1$ . It should be noted that the average overlap of a random conformation with the corresponding reference state is for the sequences considered close to  $\langle Q \rangle \approx 0.66$ . As a rule of thumb, it can be concluded that values  $Q < 0.8$  indicate weak or no significant similarity of a given structure with the reference conformation.

For the qualitative discussion of the folding characteristics, we consider the multicanonical histograms of energy  $E$  and angular overlap  $Q$ ,  $H_{\text{muca}}(E, Q) =$

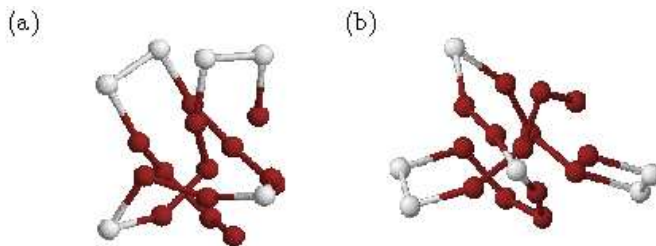
$\sum_t \delta_{E,E(\mathbf{x}_t)} \delta_{Q,Q(\mathbf{x}_t)}$ , where the sum runs over all Monte Carlo sweeps  $t$  in the multicanonical simulation, which yields a constant energy distribution  $h_{\text{muca}}(E) = \int_0^1 dQ H_{\text{muca}}(E, Q) \approx \text{const.}$  In consequence,  $H_{\text{muca}}(E, Q)$  is useful for identifying the folding channels, independently of temperature. Restricting the canonical partition function at temperature  $T$  to the “microoverlap” ensemble with overlap  $Q$ ,  $Z(Q) = \int \mathcal{D}\mathbf{X} \delta(Q - Q(\mathbf{X})) \exp\{-E(\mathbf{X})/k_B T\}$ , where the integral is over all possible conformations  $\mathbf{X}$ , we define the overlap free energy as  $F(Q) = -k_B T \ln Z(Q)$ .

Figures 15(a)–(c) show the thus obtained multicanonical histograms  $H_{\text{muca}}(E, Q)$  (left) and the overlap free-energy landscapes  $F(Q)$  (right) at different temperatures for the three sequences listed in Table 4. The different branches of  $H_{\text{muca}}(E, Q)$  indicate the channels the heteropolymer can follow in the folding process towards the reference structure. The heteropolymers, whose sequences differ only by permutations, exhibit noticeable differences in the folding behavior towards the native conformations. The first interesting observation is that the minimalistic model used is capable of revealing the different folding behaviors of the wild-type and permuted sequences. The second remarkable result is that the angular overlap parameter  $Q$  is a surprisingly manifest measure for the peptide macrostate.

From Fig. 15(a) we conclude that folding of sequence S1 exhibits a typical two-state characteristics. Above the transition temperature, i.e., in the regime of denatured conformations D, conformations possess a random-coil-like overlap  $Q \approx 0.7$ , i.e., there is no significant similarity with the reference structure. Close to  $T \approx 0.1$  the global minimum of the corresponding overlap free energy  $F(Q)$  changes discontinuously towards larger  $Q$  values, and at the transition state the denatured (D) and the folded macrostates (N) are equally probable. The existence of this pronounced transition state is a characteristic indication for first-order-like two-state folding. Decreasing the temperature further, the native-fold-like conformations ( $Q > 0.95$ ) dominate and fold smoothly towards the  $Q = 1$  reference structure, i.e., the lowest-energy conformation found for sequence S1.

The folding behavior of sequence S2 is significantly different, as Fig. 15(b) shows, and is a typical example for a folding event through an intermediate macrostate. The main channel D bifurcates and a side channel I branches off continuously. For smaller energies (or lower temperatures), this branching is followed by the formation of a third channel N, which ends in the native fold. The characteristics of folding-through-intermediates is also reflected by the free-energy landscapes. Starting at high temperatures in the pseudophase of denatured conformations D with  $Q \approx 0.76$ , the intermediary phase I with  $Q \approx 0.9$  is reached close to the temperature  $T \approx 0.05$ . Decreasing the temperature further below the native-folding threshold close to  $T = 0.01$ , the hydrophobic-core formation is finished and stable native-fold-like conformations with  $Q > 0.97$  dominate in regime N.

The most extreme behavior of the three exemplified sequences is found for sequence S3, where the main channel D does not decay in favor of a sin-



**Fig. 16.** Lowest-energy conformations for sequence S3, considered as (a) reference structure  $\mathbf{X}^{(0)}$  and (b) alternative metastable conformation, whose angular overlap with  $\mathbf{X}^{(0)}$  is  $Q \approx 0.75$ .

gle native-fold channel. In fact, in Fig. 15(c) we observe both, *two* separate native-fold channels,  $M_1$  and  $M_2$ , and a bifurcating main channel. Above the folding transition ( $T = 0.2$ ), the typical sequence-independent denatured conformations in D ( $Q \approx 0.77$ ) dominate. Annealing below the glass-transition threshold, several channels form and coexist. The two most prominent channels (to which the lowest-energy conformations belong that we found in the simulations) eventually lead for  $T \approx 0.01$  to ensembles of states  $M_1$  with  $Q > 0.97$ , which are similar to the reference structure shown in Fig. 16(a), and conformations  $M_2$  with  $Q \approx 0.75$ . The lowest-energy conformation found in this regime is shown in Fig. 16(b). It is structurally different but energetically almost degenerate compared with the reference structure. It should also be noted that the lowest-energy main-channel conformations have only slightly larger energies than the two native folds. Thus, the folding of this heteropolymer is accompanied by a very complex, amorphous folding characteristics. In fact, the multiple-peaked distribution  $H_{\text{muca}}(E, Q)$  near minimum energies is a strong indication for metastability. A native fold in the natural sense does not exist, the  $Q = 1$  conformation is only a reference structure but the folding towards this structure is not distinguished as it is in the folding characteristics of sequences S1 and S2. These results demonstrate that it is possible to find clear indications for three different folding characteristics known from real proteins by analyzing macrostates based on an angular overlap parameter within a minimalistic heteropolymer frame. The physical objective is not only on establishing a quantitative one-to-one correspondence between model and real peptides (which, in general, is not in the focus of minimalistic, effective models), but also on a more comprehensive, qualitative understanding of universal aspects of protein folding. We find that for selected hydrophobic-polar heteropolymer sequences characteristic folding behaviors such as two-state folding, folding through intermediates, and metastability can be observed which are qualitatively comparable with real folding events in nature. Beyond the general interest in understanding complex aspects of protein folding, the preparation of synthetic peptide macrostates in future applications, e.g., in the design of substrate- or pattern-selective polymers

is strongly connected with the understanding of such conformational folding transitions.

## 5 Peptide aggregation

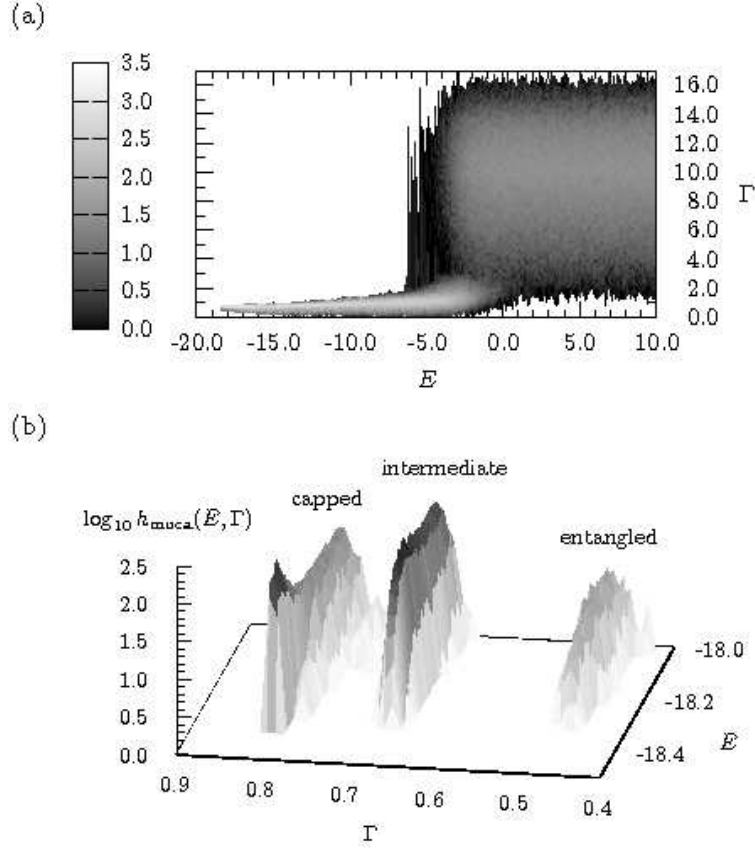
Another important, because biologically relevant, example for cooperative structure formation processes is the aggregation of proteins. A prominent example, where this process has disastrous effects, is the oligomerization of the A $\beta$  protein which is associated with Alzheimer's disease.

A mesoscopic model for the aggregation of multiple chains can simply be defined by assuming that the same type-dependent Lennard-Jones like potentials used in the single-chain form (18) describe also the inter-monomeric interaction, i.e., the interaction among monomers of different chains [86]. For the analysis of the aggregation transition let us consider the example of a complex of two identical AB peptides with sequence  $AB_2AB_2ABAB_2AB$  [86,87]. We suppose that the aggregation of the peptides should be signaled by strong fluctuations of the relative distance of the centers of masses of the individual chains. Thus we define for systems consisting of  $M$  peptides

$$I^2 = \frac{1}{2M^2} \sum_{\mu, \nu=1}^M \left( \mathbf{r}_{\text{COM}}^{(\mu)} - \mathbf{r}_{\text{COM}}^{(\nu)} \right)^2, \quad (24)$$

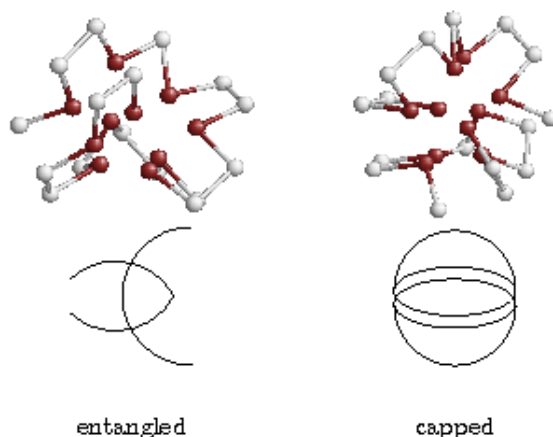
where  $\mathbf{r}_{\text{COM}}^{(\mu)}$  is the center of mass of the  $\mu$ th chain (in our example  $M = 2$ ). Actually, a multicanonical computer simulation reveals very clear indications for a single conformational transition, the aggregation transition [86,87]. This means that the peptide-peptide aggregation and the folding into a compact peptide complex are not separate transitions (at least in this example). This is illustrated in Fig. 17(a), where the color-coded multicanonical histogram as a function of energy  $E$  and the aggregation parameter  $I$  is shown. Qualitatively, two separate main branches (which are “channels” in the corresponding free-energy landscape) are apparent, between which a noticeable transition occurs. In the vicinity of the energy  $E_{\text{sep}} \approx -3.15$ , both channels overlap, i.e., the associated macrostates coexist. Since  $I$  is an effective measure for the spatial distance between the two peptides, it is obvious that conformations with separated or fragmented peptides belong to the dominating channel in the regime of high energies and large  $I$  values, whereas the aggregates are accumulated in the narrow low-energy and small- $I$  channel. Thus, the main observation from the multicanonical, comprising point of view is that the aggregation transition is a phase separation process which already appears, even for this small system, in a surprisingly clear fashion.

The high precision of the multicanonical method allows us even to see further details in the lowest-energy aggregation regime, which is usually a notoriously difficult sampling problem. Fig. 17(b) shows that the tight aggregation channel splits into three separate, almost degenerate subchannels at



**Fig. 17.** (a) Multicanonical histogram  $\log_{10} h_{\text{muca}}$  as a function of energy  $E$  and aggregation parameter  $\Gamma$ , (b) section of  $\log_{10} h_{\text{muca}}$  in the low-energy tail [87].

lowest energies. From the analysis of the conformations in this region, one finds that representative conformations with smallest  $\Gamma$  values,  $\Gamma \approx 0.45$ , are typically entangled, while those with  $\Gamma \approx 0.8$  have a spherically-capped shape. Examples are shown in Fig. 18. The also highly compact conformations belonging to the intermediate subphase do not exhibit such characteristic features and are rather globules without noticeable internal symmetries. In all cases, the aggregates contain a single compact core of hydrophobic residues. This also confirms that the aggregation is not a simple docking process of two prefolded peptides, but a complex cooperative folding-binding process. The general aggregation behavior is similar also for larger systems of more peptides with the same sequence [87].



**Fig. 18.** Representatives and schematic characteristics of entangled and spherically-capped conformations dominating the lowest-energy branches in the multicanonical histogram shown in Fig. 17(b). Dark spheres correspond to hydrophobic ( $A$ ), light ones to polar ( $B$ ) residues.

## 6 Summary

For the qualitative analysis of phase transitions, it is often sufficient to perform statistical studies of simplified effective models, where the natural complexity of the realistic system is broken down to the essential, irreducible level of cooperative behavior. The probably most famous example is the Ising model of ferromagnetism. In this model, a local short-range spin-spin interaction – which in essence is a consequence of the quantum mechanical exchange mechanism between magnetic moments – triggers in two and more dimensions a nontrivial second-order phase transition between the ordered ferromagnetic macrostate and the disordered, random paramagnetic phase. The generalization of the description of phase transitions is highly successfully achieved within the framework of Ginzburg-Landau theories, which are not only restricted to transitions of second order, but also allow investigations of symmetry breaking typically forcing first-order phase transitions. In any case, the idea is to introduce collective coordinates, or, more specific, order parameters that allow for a unique identification of the actual macrostate of the system.

The characterization of conformational (structural) transitions during the folding process of proteins is more involved as no general theory of phase transitions for finite systems is available. In fact, the finiteness of the amino acid sequence length contradicts the demand of a thermodynamic limit, which is the essential condition for thermodynamic phase transitions to occur. Nonetheless, there is hope that following a similar strategy as in the theory for phase transitions, a classification of characteristic tertiary folding transitions (e.g., single-exponential folding, two-state folding, folding through weakly stable in-

termediary states, metastability) is possible. If so, then it should be possible to construct simple models at a raw, coarse-grained level that allow firstly the introduction of unique conformational (“order”) parameters and secondly to qualitatively reproduce the known folding characteristics of classes of proteins.

As the Ising model will not be an adequate model for precise questions regarding a *specific* ferromagnet, it is also not expected that a simple, coarse-grained model will reveal the folding behavior of a *specific* protein. This means, for explaining the folding characteristics of a *specific* protein, doubtlessly a microscopic all-atom model incorporating interactions acting over all length and energy scales is required.

In this lecture, we have demonstrated, however, that results obtained from simple lattice and off-lattice heteropolymer models are indeed capable of revealing characteristic features of proteins (stability of designing sequences, designable conformations) and protein folding (folding channels, free-energy landscapes). As far as important qualitative features of peptides and proteins on intermediate length scales are concerned, such models are thus of comparable significance as the more detailed atomic descriptions.

## Acknowledgments

We are grateful to Peter Grassberger and Hsiao-Ping Hsu for detailed informations about PERM and its improved variants. We are also indebted to Anders Irbäck, Sandipan Mohanty, and Simon Mitternacht for discussions on coarse-grained and simplified variants of microscopic protein models, and to Bernd A. Berg, Ulrich H. E. Hansmann, Yuko Okamoto, Tarık Çelik, and Gökhan Gököğlu for discussions on general aspects of protein folding. We thank Handan Arkin, Reinhard Schiemann, Thomas Vogel, Stefan Schnabel, Anna Kallias, Jakob Schluttig, and Christoph Junghans for cooperation in studies of coarse-grained lattice and off-lattice heteropolymer models. The investigations of hybrid interfaces were inspired by experiments of Karsten Goede and Marius Grundmann, which we would like to thank for close collaboration. This work is partially supported by the DFG (German Science Foundation) grant under contract No. JA 483/24-1 and by the DAAD-STINT Personnel Exchange Programme with Sweden. M.B. thanks the DFG and the Wenner-Gren Foundation for research fellowships. Support by the JUMP supercomputer time grant No. h1z11 of the John von Neumann Institute for Computing (NIC), Forschungszentrum Jülich, is also gratefully acknowledged.

## References

1. T. E. Creighton, *Proteins: Structure and Molecular Properties*, 2nd. ed. (Freeman, New York, 1993).



2. C. Branden and J. Tooze, *Introduction to Protein Structure*, 2nd. ed. (Garland, New York, 1999).
3. K. A. Dill, *Protein Science* **8**, 1166 (1999).
4. C. Tang, *Physica A* **288**, 31 (2000).
5. See also the contributions of U. H. E. Hansmann, A. Irback, and Y. Okamoto to simulations of all-atom protein models in this volume.
6. B. L. de Groot, T. Frigato, V. Helms, and H. Grubmüller, *J. Mol. Biol.* **333**, 279 (2003).
7. R. A. Böckmann and H. Grubmüller, *Biophys. J.* **85**, 1482 (2003).
8. U. H. E. Hansmann, *Physica A* **254**, 15 (1998); T. Nagasima, Y. Sugita, A. Mitsutake, and Y. Okamoto, *Comp. Phys. Comm.* **146**, 69 (2002).
9. K. A. Dill, *Biochemistry* **24**, 1501 (1985); K. F. Lau and K. A. Dill, *Macromolecules* **22**, 3986 (1989).
10. B. Berger and T. Leighton, *J. Comp. Biol.* **5**, 27 (1998); P. Crescenzi, D. Goldman, C. Papadimitriou, A. Piccolboni, and M. Yannakakis, *ibid.*, 423 (1998).
11. A. Irback and E. Sandelin, *J. Chem. Phys.* **108**, 2245 (1998).
12. A. Irback and C. Troein, *J. Biol. Phys.* **28**, 1 (2002).
13. H. Cejtin, J. Edler, A. Gottlieb, R. Helling, H. Li, J. Philbin, C. Tang, and N. Wingreen, *J. Chem. Phys.* **116**, 352 (2002).
14. R. Schiemann, M. Bachmann, and W. Janke, *J. Chem. Phys.* **122**, 114705 (2005).
15. R. Schiemann, M. Bachmann, and W. Janke, *Comp. Phys. Comm.* **166**, 8 (2005).
16. K. Yue and K. A. Dill, *Phys. Rev. E* **48**, 2267 (1993); *Proc. Natl. Acad. Sci. U.S.A.* **92**, 146 (1995).
17. T. C. Beutler and K. A. Dill, *Prot. Sci.* **5**, 2037 (1996).
18. R. Unger and J. Moult, *J. Mol. Biol.* **231**, 75 (1993).
19. N. Krasnogor, W. E. Hart, J. Smith, and D. A. Pelta, *Proc. Genetic and Evolutionary Computation Conf. (GECCO99)*, Orlando, 1999, p. 1596.
20. Y. Cui, W. H. Wong, E. Bornberg-Bauer, and H. S. Chan, *Proc. Natl. Acad. Sci. U.S.A.* **99**, 809 (2002).
21. N. Lesh, M. Mitzenmacher, and S. Whitesides, *Int. Conf. Research in Computational Molecular Biology (RECOMB'03)*, Berlin, 2003, p. 188.
22. T. Jiang, Q. Cui, G. Shi, and S. Ma, *J. Chem. Phys.* **119**, 4592 (2003).
23. F. Seno, M. Vendruscolo, A. Maritan, and J. R. Banavar, *Phys. Rev. Lett.* **77**, 1901 (1996).
24. R. Ramakrishnan, B. Ramachandran, and J. F. Pekny, *J. Chem. Phys.* **106**, 2418 (1997).
25. A. Irback, C. Peterson, F. Potthast, and E. Sandelin, *Phys. Rev. E* **58**, R5249 (1998).
26. L. W. Lee and J.-S. Wang, *Phys. Rev. E* **64**, 056112 (2001).
27. G. Chikenji, M. Kikuchi, and Y. Iba, *Phys. Rev. Lett.* **83**, 1886 (1999); and references therein.
28. M. N. Rosenbluth and A. W. Rosenbluth, *J. Chem. Phys.* **23**, 356 (1955).
29. D. Aldous and U. Vazirani, *'Go with the Winners' Algorithms*, 35th Annual Symposium on Foundations of Computer Science, Santa Fe, 1994, p. 492.
30. P. Grassberger, *Phys. Rev. E* **56**, 3682 (1997).
31. H. Frauenkron, U. Bastolla, E. Gerstner, P. Grassberger, and W. Nadler, *Phys. Rev. Lett.* **80**, 3149 (1998); U. Bastolla, H. Frauenkron, E. Gerstner, P. Grassberger, and W. Nadler, *Proteins* **32**, 52 (1998).

32. P. Grassberger and W. Nadler, '*Go with the Winners' Simulations*, in: *Computational Statistical Physics – From Billiards to Monte Carlo*, edited by K. H. Hoffmann and M. Schreiber (Springer, Berlin, 2002), p. 169, and references therein.
33. H.-P. Hsu, V. Mehra, W. Nadler, and P. Grassberger, J. Chem. Phys. **118**, 444 (2003); Phys. Rev. E **68**, 21113 (2003).
34. M. Bachmann and W. Janke, Phys. Rev. Lett. **91**, 208105 (2003).
35. M. Bachmann and W. Janke, J. Chem. Phys. **120**, 6779 (2004).
36. R. J. Najmanovich, J. L. deLyra, and V. B. Henriques, Physica A **249**, 374 (1998).
37. K. Yue, K. M. Fiebig, P. D. Thomas, H. S. Chan, E. I. Shakhnovich, and K. A. Dill, Proc. Natl. Acad. Sci. U.S.A. **92**, 325 (1995).
38. E. E. Lattman, K. M. Fiebig, and K. A. Dill, Biochemistry **33**, 6158 (1994).
39. L. Toma and S. Toma, Prot. Sci. **5**, 147 (1996).
40. S. Miyazawa and R. L. Jernigan, J. Mol. Biol. **256**, 623 (1996).
41. Two sequences are only distinguished, if they are not symmetric under reversal of their residues. For a chain with length  $N = 4$ , for example, there are 10 relevant sequences instead of  $2^4 = 16$ .
42. M. Vendruscolo and E. Domany, Folding & Design **2**, 295 (1997); *ibid.* **3**, 329 (1998).
43. E. G. Emberly, J. Miller, C. Zeng, N. S. Wingreen, and C. Tang, Proteins **47**, 295 (2002).
44. H. Li, R. Helling, C. Tang, and N. Wingreen, Science **273**, 666 (1996).
45. N. Madras and A. D. Sokal, J. Stat. Phys. **50**, 109 (1988).
46. A. M. Ferrenberg and R. H. Swendsen, Phys. Rev. Lett. **63**, 1195 (1989).
47. J. C. Guillou and J. Zinn-Justin, Phys. Rev. Lett. **39**, 95 (1977); Phys. Rev. B **21**, 3976 (1980); A. Pelissetto and E. Vicari, Phys. Rep. **368**, 549 (2002).
48. B. A. Berg and T. Neuhaus, Phys. Lett. B **267**, 249 (1991), Phys. Rev. Lett. **68**, 9 (1992).
49. W. Janke, Physica A **254**, 164 (1998); B. A. Berg, Fields Inst. Comm. **26**, 1 (2000).
50. T. Vrbová and S. G. Whittington, J. Phys. A **29**, 6253 (1996); J. Phys. A **31**, 3989 (1998); T. Vrbová and K. Procházka, J. Phys. A **32**, 5469 (1999).
51. M. D. Yoder, N. T. Keen, and F. Jurnak, Science **260**, 1503 (1993).
52. M. Rief, H. Clausen-Schaumann, and H. Gaub, Nature Struct. Biol. **6**, 346 (1999).
53. D. E. Smith, S. Tans, S. Smith, S. Grimes, D. L. Anderson, and C. Bustamante, Nature **413**, 748 (2001).
54. J. J. Gray, Curr. Opin. Struct. Biol. **14**, 110 (2004).
55. E. Nakata, T. Nagase, S. Shinkai, and I. Hamachi, J. Am. Chem. Soc. **126**, 490 (2004).
56. E. Balog, T. Becker, M. Oetl, R. Lechner, R. Daniel, J. Finney, and J. C. Smith, Phys. Rev. Lett. **93**, 028103 (2004); M. Ikeguchi, J. Ueno, M. Sato, and A. Kidera, Phys. Rev. Lett. **94**, 078102 (2005).
57. J. Forsman and C. E. Woodward, Phys. Rev. Lett. **94**, 118301 (2005); G. Reiter, Phys. Rev. Lett. **87**, 186101 (2001).
58. S. Metzger, M. Müller, K. Binder, and J. Baschnagel, J. Chem. Phys. **118**, 8489 (2003).
59. T. Bogner, A. Degenhard, and F. Schmid, Phys. Rev. Lett. **93**, 268108 (2004).

60. G. M. Foo and R. B. Pandey, Phys. Rev. Lett. **80**, 3767 (1998); Phys. Rev. E **61**, 1793 (2000).
61. R. Hegger and P. Grassberger, J. Phys. A **27**, 4069 (1994).
62. Y. Singh, D. Giri, and S. Kumar, J. Phys. A **34**, L67 (2001); R. Rajesh, D. Dhar, D. Giri, S. Kumar, and Y. Singh, Phys. Rev. E **65**, 056124 (2002).
63. M. S. Causo, J. Chem. Phys. **117**, 6789 (2002).
64. J. Krawczyk, T. Prellberg, A. L. Owczarek, and A. Rechnitzer, Europhys. Lett. **70**, 726 (2005).
65. J.-H. Huang and S.-J. Han, J. Zhejiang Univ. Sci. **5**, 699 (2004).
66. M. Bachmann and W. Janke, Phys. Rev. Lett. **95**, 058102 (2005).
67. M. Bachmann and W. Janke, Phys. Rev. E **73**, 041802 (2006).
68. M. Bachmann and W. Janke, Phys. Rev. E **73**, 020901(R) (2006).
69. M. Bachmann and W. Janke, *Chain-growth simulations of lattice-peptide adsorption to attractive substrates*, in: Proceedings of the NIC Symposium 2006, John von Neumann Institute for Computing, Jülich, NIC Series vol. **32**, ed. by G. Münster, D. Wolf, and M. Kremer (NIC, Jülich, 2006), p. 245.
70. F. Celestini, T. Frisch, and X. Oyharcabal, Phys. Rev. E **70**, 012801 (2004).
71. J. Krawczyk, T. Prellberg, A. L. Owczarek, and A. Rechnitzer, J. Stat. Mech. (2004) P10004.
72. P. Benetatos and E. Frey, Phys. Rev. E **70**, 051806 (2004).
73. M. Breidenreich, R. R. Netz, and R. Lipowsky, Europhys. Lett. **49**, 431 (2000); Eur. Phys. J. E **5**, 403 (2001).
74. S. Brown, Nature Biotechn. **15**, 269 (1997).
75. R. Braun, M. Sarikaya, and K. Schulten, J. Biomater. Sci. Polym. Ed. **13**, 747 (2002).
76. S. R. Whaley, D. S. English, E. L. Hu, P. F. Barbara, A. M. Belcher, Nature (London) **405**, 665 (2000).
77. K. Goede, P. Busch, and M. Grundmann, Nano Lett. **4**, 2115 (2004).
78. N. Gupta and A. Irback, J. Chem. Phys. **120**, 3983 (2004).
79. F. H. Stillinger, T. Head-Gordon, and C. L. Hirshfeld, Phys. Rev. E **48**, 1469 (1993); F. H. Stillinger and T. Head-Gordon, Phys. Rev. E **52**, 2872 (1995).
80. R. Du, V. S. Pande, A. Yu. Grosberg, T. Tanaka, and E. S. Shakhnovich, J. Chem. Phys. **108**, 334 (1998).
81. V. S. Pande and D. S. Rokhsar, Proc. Natl. Acad. Sci. (USA) **96**, 1273 (1999).
82. U. H. E. Hansmann, M. Masuya, and Y. Okamoto, Proc. Natl. Acad. Sci. (USA) **94**, 10652 (1997).
83. B. A. Berg, H. Noguchi, and Y. Okamoto, Phys. Rev. E **68**, 036126 (2003).
84. M. Bachmann, H. Arkin, and W. Janke, Phys. Rev. E **71**, 031906 (2005).
85. S. Schnabel, M. Bachmann, and W. Janke, Phys. Rev. Lett. **98**, 048103 (2007); J. Chem. Phys. **126**, 105102 (2007).
86. C. Junghans, M. Bachmann, and W. Janke, Phys. Rev. Lett. **97**, 218103 (2006).
87. C. Junghans, M. Bachmann, and W. Janke, preprint (2007).